

# CSS Minification via Constraint Solving

(Technical Report)

Matthew Hague<sup>1</sup>, Anthony W. Lin<sup>2,3</sup> and Chih-Duo Hong<sup>3</sup>

<sup>1</sup> Royal Holloway, University of London, UK

<sup>2</sup> Technische Universität Kaiserslautern, Germany

<sup>3</sup> University of Oxford, UK

## Abstract

Minification is a widely-accepted technique which aims at reducing the size of the code transmitted over the web. This paper concerns the problem of semantic-preserving minification of Cascading Style Sheets (CSS) — the de facto language for styling web documents — based on merging similar rules.

The cascading nature of CSS makes the semantics of CSS files sensitive to the ordering of rules in the file. To automatically identify rule-merging opportunities that best minimise file size, we reduce the rule-merging problem to a problem concerning “CSS-graphs”, i.e., node-weighted bipartite graphs with a dependency ordering on the edges, where weights capture the number of characters.

Constraint solving plays a key role in our approach. Transforming a CSS file into a CSS-graph problem requires us to extract the dependency ordering on the edges (an NP-hard problem), which requires us to solve the selector intersection problem. To this end, we provide the first full formalisation of CSS3 selectors (the most stable version of CSS) and reduce their selector intersection problem to satisfiability of quantifier-free integer linear arithmetic, for which highly-optimised SMT-solvers are available. To solve the above NP-hard graph optimisation problem, we show how Max-SAT solvers can be effectively employed. We have implemented our rule-merging algorithm, and tested it against approximately 70 real-world examples (including examples from each of the top 20 most popular websites). We also used our benchmarks to compare our tool against six well-known minifiers (which implement other optimisations). Our experiments suggest that our tool produced larger savings. A substantially better minification rate was shown when our tool is used together with these minifiers.

## 1 Introduction

*Minification* [62, Chapter 12] is a widely-accepted technique in the web programming literature that aims at decreasing the size of the code transmitted over the web, which can directly improve the response-time performance of a website. Page load time is of course crucial for users’ experience, which impacts the performance of an online business and is increasingly being included as a ranking factor by search engines [61]. Minification bears a resemblance to traditional code compilation. In particular, it is applied only *once* right before deploying the website (therefore, its computation time does *not* impact the page load time). However, they differ in at least two ways. First, the source and target languages for minification are the same (high-level) languages. The code to which minification can be applied is typically JavaScript or CSS, but it can also be HTML, XML, SVG, etc. Second, minification applies various semantic-preserving transformations with the objective of reducing the size of the code.

This paper concerns the problem of minifying CSS (Cascading Style Sheets), which is the de facto language for styling web documents (HTML, XML, etc.) as developed and maintained by the World Wide Web Consortium (W3C) [4]. We will minify the CSS without reference to the documents it may be designed to style. We refer the reader to Section 9 for a discussion of document-independent and document-dependent related work.

A CSS file consists of a list of CSS *rules*, each containing a list of *selectors* — each selecting nodes in the *Document Object Model* (DOM), which is a tree structure representing the document — and a list of *declarations*, each assigning values to selected nodes' display attributes (e.g. `blue` to the property `color`). Real-world CSS files can easily have many rules (in the order of magnitude of 1000), e.g., see the statistics for popular sites [1] on <http://cssstats.com/>. As Souders wrote in his book [62], which is widely regarded as the bible of web performance optimisation:

The greatest potential for [CSS file] size savings comes from optimizing CSS — merging identical classes, removing unused classes, etc. This is a complex problem, given the order-dependent nature of CSS (the essence of why it's called *cascading*). This area warrants further research and tool development.

More and more CSS minifiers have been, and are being, developed. These include `YUI Compressor` [59], `cssnano` [13], `minify` [16], `clean-css` [51], `csso` [19], and `cssmin` [8], to name a few. Such CSS minifiers apply various syntactic transformations, typically including removing whitespace characters and comments, and using abbreviations (e.g. `#f60` instead of `#ff6600` for the colour orange).

In this paper, we propose a new class of CSS transformations based on merging similar or duplicate rules in the CSS file (thereby removing repetitions across multiple rules in the file) that could reduce file size while preserving the rendering information. To illustrate the type of transformations that we focus on in this paper (a more detailed introduction is given in Section 2), consider the following simple CSS file.

```
#a { color:red; font-size:large }
#c { color:green }
#b { color:red; font-size:large }
```

The selector `#a` selects a node with *ID* `a` (if it exists). Note that the first and the third rules above have the same property declarations. Since one *ID* can be associated with at most one node in a DOM-tree, we can merge these rules resulting in the following equivalent file.

```
#a, #b { color:red; font-size:large }
#c { color:green }
```

Identifying such a rule-merging-based minification opportunity — which we shall call *merging opportunity* for short — in general is non-trivial since a CSS file is sensitive to the ordering of rules that may match the same node, i.e., the problem mentioned in Souders' quote above. For example, let us assume that the three selectors `#a`, `#b`, and `#c` in our CSS example instead were `.a`, `.b`, and `.c`, respectively. The selector `.a` selects nodes with *class* `a`. Since a class can be associated with multiple nodes in a DOM-tree, the above merging opportunity could change how a page is displayed and therefore would not be valid, e.g., consider a page with an element that has two classes, `.b` and `.c`. This element would be displayed as red by the original file (since red appears later), but as green by the file after applying the transformation. Despite this, in this case we would still be able to merge the *subrules* of the first and the third rules (associated with the `font-size` property) resulting in the following smaller equivalent file:

```
.a { color:red }
.c { color:green }
.b { color:red }
.a, .b { font-size:large }
```

This example suggests two important issues. First, identifying a merging opportunity in a general way requires a means of checking whether two given selectors may *intersect*, i.e., select the same node in *some* document tree. Second, a given CSS file could have a large number of merging opportunities at the same time (the above example has at least two). Which one shall we pick? There are multiple answers to this question. Note first that merging rules can be iterated multiple times to obtain increasingly smaller CSS files. In fact, we can define an optimisation problem which, given a CSS file, computes a sequence of applications of semantic-preserving merging rules that yields a *globally minimal* CSS file. Finding a provably globally minimal CSS file is a computationally difficult problem that seems well beyond the reach of current technologies (constraint solving, and otherwise). Therefore, we propose to use the simple *greedy strategy*: while there is a merging opportunity that can make the CSS file smaller, apply an *optimal* one, i.e., a merging opportunity that reduces the file size the most<sup>1</sup>. There are now two problems to address. First, can we efficiently find an optimal merging opportunity for a given CSS file? This paper provides a positive answer to this question by exploiting a state-of-the-art constraint solving technology. Second, there is a potential issue of getting stuck at a *local minimum* (as is common with any optimisation method based on gradient descent). Does the greedy strategy produce a meaningful space saving? As we will see in this paper, the greedy approach could already produce space savings that are beyond the reach of current CSS minification technologies.

## 1.1 Contributions

We first formulate a *general class of semantic-preserving transformations* on CSS files that captures the above notion of merging “similar” rules in a CSS file. Such a program transformation has a clean graph-theoretic formulation (see Section 4). Loosely speaking, a CSS rule corresponds to a biclique (complete bipartite graph)  $B$ , whose edges connect nodes representing selectors and nodes representing property declarations. Therefore, a CSS file  $F$  corresponds to a sequence of bicliques that covers all of the edges in the bipartite graph  $\mathcal{G}$  which is constructed by taking the (non-disjoint) union of all bicliques in  $F$ . Due to the cascading nature of CSS, the file  $F$  also gives rise to an (implicit) ordering  $\prec$  on the edges of  $\mathcal{G}$ . Therefore, any new CSS file  $F'$  that we produce from  $F$  must also be a valid covering of the edges of  $\mathcal{G}$  and respect the order  $\prec$ . As we will see, the above notion of merging opportunity can be defined as a pair  $(\bar{B}, j)$  of a new rule  $\bar{B}$  and a position  $j$  in the file, and that applying this transformation entails inserting  $\bar{B}$  in position  $j$  and removing all redundant nodes (i.e. either a selector or a property declaration) in rules at position  $i < j$ .

Several questions remain. First is how to compute the edge order  $\prec$ . The core difficulty of this problem is to determine whether two CSS selectors can be matched by the same node in some document tree (a.k.a. the *selector intersection problem*). Second, among the multiple potential merging opportunities, how do we automatically compute a rule-merging opportunity that best minimises the size of the CSS file. We provide solutions to these questions in this paper.

**Computing the edge order  $\prec$ .** In order to handle the selector intersection problem, we first provide a complete formalisation of CSS3 selectors [15] (currently the most stable version of CSS selectors). We then give a polynomial-time reduction from the selector intersection problem to satisfiability over quantifier-free theory of integer linear arithmetic, for which highly optimised SMT-solvers (e.g. Z3 [18]) are available. To achieve this reduction, we provide a chain of polynomial-time reductions. First, we develop a new class of *automata over data trees* [9], called *CSS automata*, which can capture the expressivity of CSS selectors. This reduces the selector intersection problem to the language intersection problem for CSS automata. Second, unlike the case for CSS selectors, the languages recognised by

<sup>1</sup>Not to be confused with the smallest CSS file that is equivalent with the original file, which in general cannot be obtained by applying a single merging

CSS automata enjoy closure under intersection. [.b .a and .c .a are individually CSS selectors, however their conjunction is not a CSS selector<sup>2</sup>.] This reduces language intersection of CSS automata to language non-emptiness of CSS automata. To finish this chain of reductions, we provide a reduction from the problem of language non-emptiness of CSS automata to satisfiability over quantifier-free theory of integer linear arithmetic. The last reduction is technically involved and requires insights from logic and automata theory, which include several small model lemmas (e.g. the sufficiency of considering trees with a small number of node labels that can be succinctly encoded). This is despite the fact that CSS selectors may force the smallest satisfying tree to be exponentially big.

**Formulation of the “best” rule-merging opportunity and how to find it.** Since our objective is to minimise file size, we may equip the graph  $\mathcal{G}$  by a *weight function*  $w_t$  which maps each node to the number of characters used to define it (recall that a node is either a selector or a property declaration). Hence, the function  $w_t$  allows us to define the size of a CSS file  $F$  (i.e. a covering of  $\mathcal{G}$  respecting  $\prec$ ) by taking the sum of weights  $w_t(v)$  ranging over all nodes  $v$  in  $F$ . The goal, therefore, is to find a merging opportunity  $(\bar{B}, j)$  of  $F$  that produces  $F'$  with a minimum file size. We show how this problem can be formulated as a (partially weighted) Max-SAT instance [3] in such a way that several existing Max-SAT solvers (including Z3 [7] and MaxRes [46]) can handle it efficiently. This Max-SAT encoding is non-trivial: the naive encoding causes Max-SAT solvers to require prohibitively long run times even on small examples. A naive encoding would allow the Max-SAT solver to consider any rule constructed from the nodes of the CSS file, and then contain clauses that prohibit edges that do not exist in the original CSS file (as these would introduce new styling properties). Our encoding forces the Max-SAT solver to only consider rules that do not introduce new edges. We do this by enumerating all *maximal* bicliques in the graph  $\mathcal{G}$  (maximal with respect to set inclusion) and developing an encoding that allows the Max-SAT solver to only consider rules that are contained within a maximal biclique. We employ the algorithm from [36] for enumerating all maximal bicliques in a bipartite graph, which runs in time polynomial in the size of the input and output. Therefore, to make this algorithm run efficiently, the number of maximal bicliques in the graph  $\mathcal{G}$  cannot be very large. Our benchmarking (using approximately 70 real-world examples including CSS files from each of the top 20 websites [1]) suggests that this is the case for graphs  $\mathcal{G}$  generated by CSS files (with the maximum ratio between the number of bicliques and the number of rules being 2.05). Our experiments suggest that the combination of the enumeration algorithm of [36] and Z3 [7] makes the problem of finding the best merging opportunity for CSS files practically feasible.

**Experiments** We have implemented our CSS minification algorithm in the tool SATCSS which greedily searches for and applies the best merging opportunity to a given CSS file until no more rule-merging can reduce file size. This will be made available in the supplementary material. The source code and a working disk image are also available from the following URLs:

<https://github.com/matthewhague/sat-css-tool>  
<http://www.cs.rhul.ac.uk/home/hague/sat-css-tool.img.txz>

Our tool utilises Z3 [18, 7] as a backend solver for Max-SAT and SMT over integer linear arithmetic. We have tested our tool on around 70 examples from real-world websites (including examples from each of the top 20 websites [1]) with promising experimental results. We found that SATCSS (which only performs rule-merging) yields larger savings on our benchmarks in comparison to six popular CSS minifiers [60], which support many other optimisations but not rule-merging. More precisely, when run individually, SATCSS reduced the file size by a third quartile of 6.90% and a median value of

<sup>2</sup>The conjunction can be expressed by the *selector group* .b .c .a, .c .b .a, .b.c .a.

3.79%. The six mainstream minifiers achieved savings with third quantiles and medians up to 5.45% and 3.29%, respectively. Since these are orthogonal minification opportunities, one might suspect that applying these optimisations together could further improve the minification performance, which is what we discover during our experiments. More precisely, when we run our tool after running any one of these minifiers, the third quartile of the savings can be increased to 8.26% and the median to 4.70%. The additional gains obtained by our tool on top of the six minifiers (as a percentage of the original file size) have a third quartile of 5.03% and a median value of 2.80%. Moreover, the ratios of the percentage of savings made by SATCSS to the percentage of savings made by the six minifiers have third quartiles of at least 136% and medians of at least 48%. In fact, in the case of `cleancss` which has a third quartile saving of 5.26% and median saving of 3.03%, applying SATCSS thereafter results in a third quartile saving of 9.50% and median saving of 6.10%. These figures clearly indicate a substantial proportion of extra space savings made by SATCSS. See Table 1 and Figure 15 for more detailed statistics.

## 1.2 Organisation

Section 2 provides a gentle and more detailed introduction (by example) to the problem of CSS minification via rule-merging. Preliminaries (notation, and basic terminologies) can be found in Section 3. In Section 4 we formalise the rule-merging problem in terms of what we will call CSS-graphs. In Section 5, we provide a formalisation of CSS3 selectors and an outline of an algorithm for solving their intersection problem, from which we can extract the edge order of a CSS-graph that is required when reducing the rule-merging problem to the edge-covering problem of CSS-graphs. Since the algorithm solving the selector intersection problem is rather intricate, we dedicate one full section (Section 6) for it. In Section 7 we show how Max-SAT solvers can be exploited to solve the rule-merging problem of CSS-graphs. We describe our experimental results in Section 8. In Section 9 we give a detailed discussion of related work. Finally, we conclude in Section 10. Additional details can be found in the appendix. The Python code and benchmarks for our tool have been included in the supplementary material, with a brief user guide. A full artefact, including virtual machine image will be included when appropriate. These are presently available from the URLs above.

## 2 CSS Rule-Merging: Example and Outline

In this section, we provide a gentler and more detailed introduction to CSS minification via merging rules, while elaborating on the difficulties of the problem. We also give a general overview of the algorithm, which may serve as a guide to the article. We begin by giving a basic introduction to HTML and CSS. Our formal models of HTML and CSS are given in Section 5.

### 2.1 HTML and CSS

In this section we give a simple example to illustrate HTML and CSS. Note, in this article our analysis takes as input a CSS file only. We cover HTML here to aid the description of CSS.

An HTML document is given in Figure 1. The format is similar to (but not precisely) XML, and is organised into a tree structure. The root node is the `html element` (or `tag`), and its two children are the `head` and `body` nodes. These contain page header information (particularly the title), which is not directly displayed, and the main page contents respectively.

The `body` node contains two children which represent the page content. The first child is a `div` element, which is used to group together parts of the page for the convenience of the developer. In this case, the developer has chosen to collect all parts of the page containing the displayed heading of the

```

<html>
  <head><title>Shopping List</title></head>
  <body>
    <div id="heading">
      
      <h1>An Example HTML Document</h1>
    </div>
    <ul>
      <li id="apple" class="fruit">Apple</li>
      <li id="broccoli">Broccoli</li>
    </ul>
  </body>
</html>

```

Figure 1: A simple HTML document.

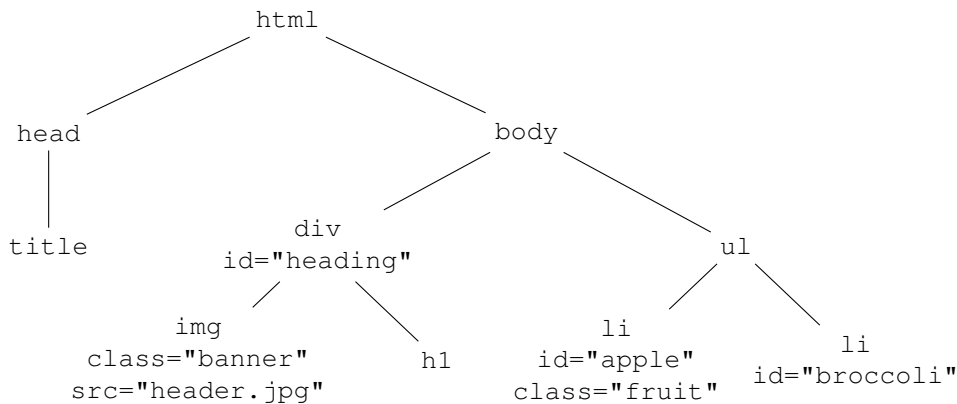


Figure 2: The HTML file in Figure 1 drawn as a tree.

page into one `div`. This `div` has an *id* which takes the value `heading`. Values specified in this way are called *attributes*. *ID* attributes uniquely identify a given node in the HTML, that is, no two nodes should have the same ID. The `div` has two children, which are an image (`img`) element and a textual heading (`h1`). The source of the image is given by `src`, which is also an attribute. The image also has a `class` attribute which takes the value `banner`. The class attribute is used to distinguish nodes which have the same tag. For example, the developer may want an image with the class `banner` to be displayed differently to an image with the class `illustration`.

The second child of the `body` node is an unordered list (`ul`). This lists contains two food items, with appropriate ID and class attributes.

To emphasize the tree structure, we give in Figure 2 the same HTML document drawn as a tree. We will give a refined version of this diagram in Section 5 when defining DOM trees formally.

A CSS file consists of a sequence of *rules*, each of the form

$$\textit{selectors} \{ \textit{declarations} \}$$

where *selectors* contains one or more (node-)selectors (separated by commas) and *declarations* contains a sequence of (visual) property declarations (separated by semicolons). An example CSS file containing two rules is given in Figure 3.

The semantics of a rule is simple: if a node can be matched by at least one of the selectors, then label the node by all the visual properties in the rule. Both rules in Figure 3 have one selector. The selector

```
img.banner { width: 100% }
#heading h1 { font-size: 30pt; font-weight: bold }
```

Figure 3: An example CSS file.

in the first rule `img.banner` matches all `img` elements which have the class `banner`. The notation `.` is used to indicate a class. Thus, this rule will match the `img` node in our example document and will specify that the image is to be wide.

The second selector demonstrates that selectors may reason about the structure of the tree. The `#heading` means that the match should begin at the node with ID `heading`. Then, the space means the match should move to any descendent of the `#heading` node, i.e., any node contained within the `div` with ID `heading` in our example page. Next, the `h1` is used to choose `h1` elements. Thus, this selector matches any node with element `h1` that occurs below a node with ID `heading`. The text of any such match should be rendered in a 30pt bold font.

There are many other features of selectors. These are described in full in Section 5.

## 2.2 CSS Rule-Merging by Example

We will now discuss more advanced features of CSS files and give an example of rule-merging. Figure 4 contains a simple CSS file (with four rules).

```
#apple { color:blue; font-size:small }
.fruit, #broccoli { color:red; font-size:large }
#orange { color:blue }
#tomato { color:red; font-size:large;
         background-color:lightblue }
```

Figure 4: A simple example of a CSS file.

The sequence of rules in a file is applied to a node  $v$  in a “cascading” fashion (hence the name Cascading Style Sheets). That is, read the rules from top to bottom and check if  $v$  can be matched by at least one selector in the rule. If so, assign the properties to  $v$  (perhaps overriding previously assigned properties, e.g., `color`) provided that the selectors matching  $v$  in the current rule have higher “values” (a.k.a. *specificities*) than the selectors previously matching  $v$ . Intuitively, the specificity of a selector [15] can be calculated by taking a weighted sum of the number of classes, IDs, tag names, etc. in the selector, where IDs have higher weights than classes, which in turn have higher weights than tag names.

For example, let us apply this CSS file to a node matching `.fruit` and `#apple`. In this case, the selectors in the first (`#apple`) and the second rules (`.fruit`) are applicable. However, since `#apple` has a higher specificity than `.fruit`, the node gets labels `color:blue` and `font-size:small`.

Two syntactically different CSS files could be “semantically equivalent” in the sense that, on each DOM tree  $T$ , both CSS files will precisely yield the same tree  $T'$  (which annotates  $T$  with visual information). In this paper, we only look at semantically equivalent CSS files that are obtained by merging similar rules. Given a CSS rule  $R$ , we say that it is *subsumed* by a CSS file  $F$  if each possible selector/property combination in  $R$  occurs in some rule in  $F$ . Notice that a rule can be subsumed in a CSS file  $F$  without occurring in  $F$  as one of the rules. For example, the rule  $R_1$

```
.fruit, #broccoli, #tomato { color:red; font-size:large }
```

is subsumed in our CSS file example (not so if `background-color:lightblue` were added to  $R_1$ ). A (*rule-merging opportunity*) consists of a CSS rule subsumed in the CSS file and a position in the

```
#apple { color:blue; font-size:small }
.fruit, #broccoli { color:red; font-size:large }
#orange { color:blue }
#tomato { color:red; font-size:large;
         background-color:lightblue }
.fruit, #broccoli, #tomato { color:red; font-size:large }
```

(a) The CSS file in Figure 4 with a rule inserted at the end.

```
#apple { color:blue; font-size:small }
#orange { color:blue }
#tomato { background-color:lightblue }
.fruit, #broccoli, #tomato { color:red; font-size:large }
```

(b) The result of trimming the CSS file in Figure 5a.

Figure 5: An example of insertion and trimming.

file to insert the rule into. An example of a merging opportunity in our CSS file example is the rule  $R_1$  and the end of the file as the insertion position. This results in a new (bigger) CSS file shown in Figure 5a.

We then “trim” this resulting CSS file by eliminating “redundant” subrules. For example, the second rule is redundant since it is completely contained in the new rule  $R_1$ . Also, notice that the subrule:

```
#tomato { color:red; font-size:large }
```

of the fourth rule in the original file is also completely redundant since it is contained in  $R_1$ . Removing these, we obtain the trimmed version of the CSS file, which is shown in Figure 5b. This new CSS file contains fewer bytes.

We may apply another round of rule-merging by inserting the rule

```
#apple, #orange { color:blue }
```

at the end of the second rule (and trim). The resulting file is even smaller. Computing such merging opportunities (i.e. which yields maximum space saving) is difficult.

Two remarks are in order. Not all merging opportunities yield a smaller CSS file. For example, if `#tomato` is replaced by the fruit vegetable

```
#vigna_unguiculata_subsp_sesquipedalis
```

the first merging opportunity above would have resulted in a larger file, which we should *not* apply. The second remark is related to the issue of *order dependency*. Suppose we add the selector `.vegetable` to the third rule in the original file, resulting in the following rule

```
.vegetable, #orange { color:blue }
```

Any node labeled by `.fruit` and `.vegetable` but no IDs will be assigned `color:blue`. However, using the first merging opportunity above yields a CSS file which assigns `color:red` to fruit vegetables, i.e., *not* equivalent to the original file. To tackle the issue of order dependency, we need to account for selectors specificity and whether two selectors with the same specificity can *intersect* (i.e. can be matched by a node in some DOM tree). Although for our examples the problem of selector intersection is quite obvious, this is not the case for CSS selectors in general. For example, the following two selectors are from a real-world CSS example found on The Guardian website.



```
.commercial--masterclasses .lineitem:nth-child(4)
.commercial--soulmates:nth-child(n+3)
```

Interested readers unfamiliar with CSS may refer to Section 5.2 for a complete definition of selectors. The key feature is the use of `:nth-child`. In the first selector it is used to only match nodes that are the fourth child of some node. For a node to match the second selector, there must be some  $n \geq 0$  such that the node is the  $(n + 3)$ th child of some node. I.e. the third, fourth, fifth, etc. These two selectors have a non-empty intersection, while changing `n+3` to `2n+3` would yield the intersection empty!

## 2.3 CSS Rule-Merging Outline

The component parts of our algorithm are given in the schematic diagram in Figure 6.

From an input CSS file, the first step is to construct a formal model of the CSS that we can manipulate. This involves extracting the selector and property declaration pairs that appear in the file. Then the edge ordering, which records the order in which selector/declaration pairs should appear, needs to be built. To compute the edge ordering, it is important to be able to test whether two selectors may match the same node in some document tree. Thus, computing the edge ordering requires an algorithm for computing whether the intersection of two selectors is empty. This process is described in detail in Section 5 and Section 6.

Once a model has been constructed, it is possible to systematically search for semantics-preserving rule-merging opportunities that can be applied to the file to reduce its overall size. The search is formulated as a MaxSAT problem and a MaxSAT solver is used to find the merging opportunity that represents the largest size saving. This encoding is described in Section 7. If a merging opportunity is found, it is applied and the search begins for another merging opportunity. If none is found, the minimised CSS file is output.

## 3 Preliminaries

### 3.1 Maths

As usual,  $\mathbb{Z}$  denotes the set of all integers. We use  $\mathbb{N}$  to denote the set  $\{0, 1, \dots\}$  of all natural numbers. Let  $\mathbb{N}_{>0} = \mathbb{N} \setminus \{0\}$  denote the set of all positive integers. For an integer  $x$  we define  $|x|$  to be the absolute value of  $x$ . For two integers  $i, j$ , we write  $[i, j]$  to denote the set  $\{i, \dots, j\}$ . Similar notations for open intervals will also be used for integers (e.g.  $(i, j)$  to mean  $\{i + 1, \dots, j - 1\}$ ). For a set  $S$ , we will write  $S^*$  (resp.  $S^+$ ) to denote the set of sequences (resp. non-empty sequences) of elements from  $S$ . When the meaning is clear, if  $S$  is a singleton  $\{s\}$ , we will denote  $\{s\}^*$  (resp.  $\{s\}^+$ ) by  $s^*$  (resp.  $s^+$ ). Given a (finite) sequence  $\sigma = s_1, \dots, s_n$ ,  $i \in [0, n]$ , and a new element  $s$ , we write  $\sigma[s \rightarrow i]$  to denote the new sequence  $s_1, \dots, s_i, s, s_{i+1}, \dots, s_n$ , i.e., inserting the element  $s$  right after the position  $i$  in  $\sigma$ .

### 3.2 Trees

We review a standard formal definition of (*rooted*) *unranked ordered trees* [22, 37, 47] from the database research community, which use it to model XML. We will use this in Section 5 to define document trees. “Unranked” means that there is no restriction on the number of children of a node, while “ordered” means that the children of each node are linearly ordered from the left-most child to the right-most child. An unranked ordered tree consists of a tree domain and a labelling, which we define below.

We identify a node by the unique path in the tree to the given node. A *tree domain* defines the set of nodes in the tree and is a set of sequences of natural numbers  $D \subseteq (\mathbb{N}_{>0})^*$ . The empty sequence is the

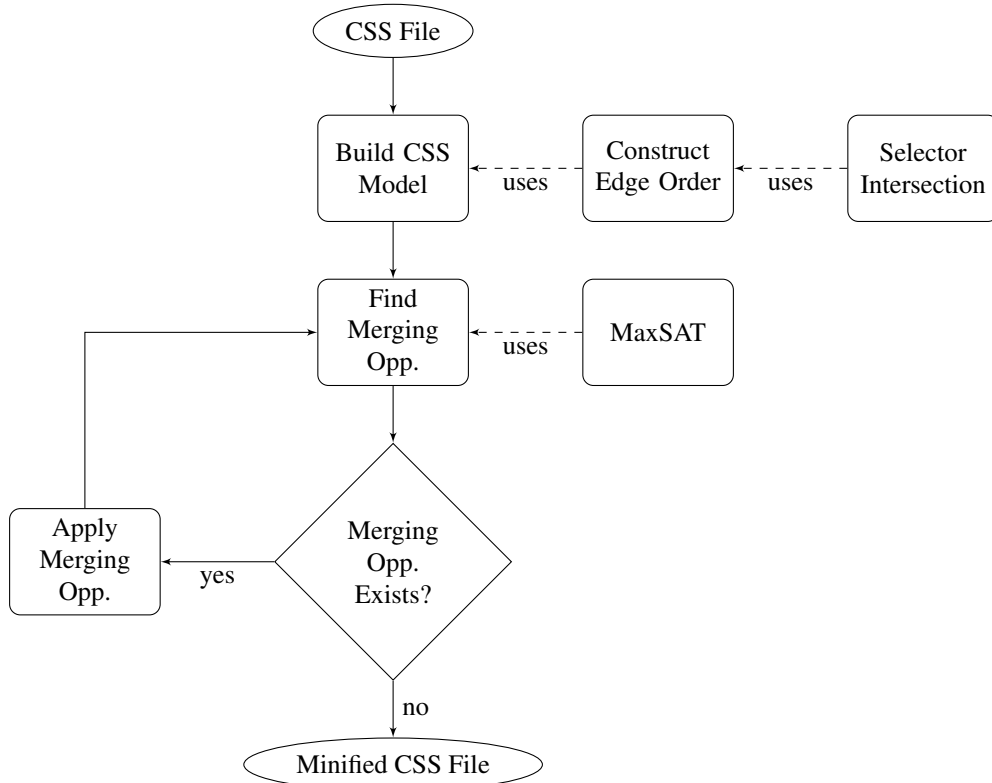


Figure 6: A Schematic Diagram of Our Approach

root node. The sequence 1 would be the first child of the root node, while 2 would be the second child. The sequence 21 would denote the first child of the second child of the root node, and so on. The tree in Figure 2 has  $D = \{\varepsilon, 1, 2, 11, 21, 22, 211, 212, 221, 222\}$ .

We require that  $D$  is both *prefix-closed* and *preceding-sibling closed*. By prefix-closed we formally mean  $\eta\nu \in D$  implies  $\eta \in D$ ; this says that the parent of each node is also a node in the tree. By preceding-sibling closed we formally mean  $\eta\nu \in D$  implies  $\eta\nu' \in D$  for all  $\nu' < \nu$ ; for example, this means that if a node has a second child, it also has a first. Observe, we write  $\eta$  for a tree node (element of  $D$ ) and  $\nu$  for an element of  $\mathbb{N}_{>0}$ .

Our trees will also be labelled by items such as element names and attribute values. In general, a  $\Sigma$ -labelled tree is a pair  $T = (D, \lambda)$  where  $D$  is a tree domain, and  $\lambda : D \rightarrow \Sigma$  is a labelling function of the nodes of  $T$  with items from a set of labels  $\Sigma$ . A simple encoding of the tree in Figure 2 will have the labelling

```

λ(ε)   = html
λ(1)   = head
λ(2)   = body
λ(11)  = title
λ(21)  = div id="heading"
λ(22)  = ul
λ(211) = img class="banner" src="header.jpg"
λ(212) = h1
λ(221) = li id="apple" class="fruit"
λ(222) = li id="broccoli".
  
```

Note, in Section 5, we will use a more complex labelling of DOM trees to fully capture all required features.

Next we recall terminologies for relationships between nodes in trees. To avoid notational clutter, we deliberately choose notation that resembles the syntax of CSS, which we define in Section 5. In the following, take  $\eta, \eta' \in D$ . We write  $\eta \otimes \eta'$  if  $\eta$  is a (strict) ancestor of  $\eta'$ , i.e., there is some  $\eta'' \in \mathbb{N}_{>0}^+$  such that  $\eta' = \eta\eta''$ . We write  $\eta \odot \eta'$  if  $\eta$  is the parent of  $\eta'$ , i.e., there is some  $\iota \in \mathbb{N}_{>0}$  such that  $\eta' = \eta\iota$ . We write  $\eta \oplus \eta'$  if  $\eta$  is the direct preceding sibling of  $\eta'$ , i.e., there is some  $\eta''$  and  $\iota \in \mathbb{N}_{>0}$  such that  $\eta = \eta''(\iota - 1)$  and  $\eta' = \eta''\iota$ . We write  $\eta \ominus \eta'$  if  $\eta$  is a preceding sibling of  $\eta'$ , i.e., there is some  $\eta''$  and  $\iota, \iota' \in \mathbb{N}_{>0}$  with  $\iota < \iota'$  such that  $\eta = \eta''\iota$  and  $\eta' = \eta''\iota'$ .

### 3.3 Max-SAT

In this paper, we will reduce the problem of identifying an optimal merging opportunity to partial weighted Max-SAT [3]. Partial weighted Max-SAT formulas are boolean formulas in CNF with *hard constraints* (a.k.a. clauses that must be satisfied) and *soft constraints* (a.k.a. clauses that may be violated with a specified cost or weight). A minimal-cost solution is the goal. Note that our clauses will not be given in CNF, but standard satisfiability-preserving conversions to CNF exist (e.g. see [12]) which straightforwardly extend to partial weighted Max-SAT.

We will present Max-SAT problems in the following form  $(\Pi_H, \Pi_S)$  where

- $\Pi_H$  are the hard constraints – that is, a set of boolean formulas that *must* be satisfied – and
- $\Pi_S$  are the soft constraints – that is a set of pairs  $(\varphi, \omega)$  where  $\varphi$  is a boolean formula and  $\omega \in \mathbb{N}$  is the weight of the constraint.

Intuitively, the weight of a soft constraint is the cost of not satisfying the constraint. The partial weighted Max-SAT problem is to find an assignment to the boolean variables that satisfies all hard constraints and minimises the sum of the weights of unsatisfied soft constraints.

### 3.4 Existential Presburger Arithmetic

In this paper, we present several encodings into *existential Presburger arithmetic*, also known as the *quantifier-free theory of integer linear arithmetic*. Here, we use extended existential Presburger formulas  $\exists x_1, \dots, x_k. \varphi$  where  $\varphi$  is a boolean combination of expressions  $\sum_{i=1}^k a_i x_i \sim b$  for constants  $a_1, \dots, a_k, b \in \mathbb{Z}$  and  $\sim \in \{\leq, \geq, <, >, =\}$  with constants represented in binary. A formula is satisfiable if there is an assignment of a non-negative integer to each variable  $x_1, \dots, x_k$  such that  $\varphi$  is satisfied. For example, a simple existential Presburger formula is shown below

$$\exists x, y, z. 0 > 2y + z - x \wedge 0 > z - y$$

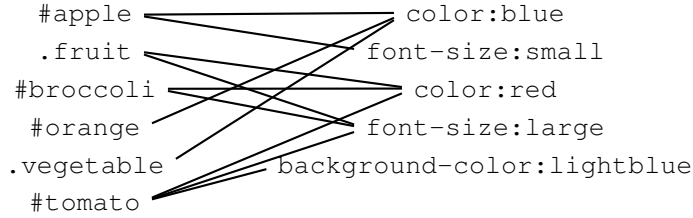
which is satisfied by any assignment to the variables  $x, y$ , and  $z$  such that  $x < 2y + z$  and  $y > z$ . The assignment  $x = 2, y = 3, z = 0$  is one such satisfying assignment. Note, when writing existential Presburger formulas, we will allow formulas that do not strictly match the format specified above. This is to aid readability and we will always be able to rewrite the formulas into the correct format. The above example formula may be written

$$\exists x, y, z. x < 2y + z \wedge y < z.$$

It is well-known that satisfiability of existential Presburger formulas is NP-complete even with the above extensions (cf. [55]). Such problems can be solved efficiently by SMT solvers such as Z3.

## 4 Formal definition of CSS rule-merging and minification

The semantics of a CSS file can be formally modelled as a CSS-graph. A *CSS-graph* is a 5-tuple  $\mathcal{G} = (S, P, E, \prec, \text{wt})$ , where  $(S, P, E)$  is a bipartite graph (i.e. the set  $V$  of vertices is partitioned into two sets  $S$  and  $P$  with  $E \subseteq S \times P$ ),  $\prec \subseteq E \times E$  gives the order dependency on the edges, and  $\text{wt} : S \cup P \rightarrow \mathbb{Z}_{>0}$  is a weight function on the set of vertices. Vertices in  $S$  are called *selectors*, whereas vertices in  $P$  are called *properties*. For example, the CSS graph corresponding to the CSS file in Figure 4 with the selector `.vegetable` added to the third rule is the following bipartite graph



such that the weight of a node is  $1 + (\text{length of the text})$ , e.g.,  $\text{wt}(\#orange) = 1 + 7 = 8$ . The reason for the extra  $+1$  is to account for the selector/property separators (i.e. commas or semi-colons), as well as the character ‘{’ (resp. ‘}’) at the end of the sequence of selectors (resp. properties). That is, in a rule, selectors are followed by a comma if another selector follows, or ‘{’ if it is the last selector, and properties are followed by a semi-color if another property follows, or ‘}’ if it is the last property declaration. We refer to  $\prec$  as the *edge order* and it intuitively states that one edge should appear strictly before the other in any CSS file represented by the graph. In this case we have  $(.fruit, \text{color:red}) \prec (.vegetable, \text{color:blue})$  because any node labeled by `.fruit` and `.vegetable` but no IDs should be assigned the property `color:blue`. There are no other orderings since each node can have at most one ID<sup>3</sup> and `.fruit` and `.vegetable` are the selectors of the lowest specificity in the file. More details on how to compute  $\prec$  from a CSS file are given in Section 5.4.

A *biclique* in  $\mathcal{G}$  is a complete bipartite subgraph, i.e., a pair  $B = (X, Y)$  of a nonempty set  $X \subseteq S$  of selectors and a nonempty set  $Y \subseteq P$  of properties such that  $X \times Y \subseteq E$  (i.e. each pair of a selector and a property in the rule is an edge). A (*CSS*) *rule* is a pair  $\bar{B} = (B, \triangleleft)$  of a biclique and a total order  $\triangleleft$  on the set of properties. The reason for the order on the properties, but not on the selectors, is illustrated by the following example of a CSS rule:

```
.a, .b { color:red; color:rgba(255,0,0,0.5) }
```

That is, nodes matching `.a` or `.b` are assigned a semi-transparent red with solid red being defined as a *fallback* when the semi-transparent red is not supported by the document reader. Therefore, swapping the order of the properties changes the semantics of the rule, but swapping the order of the selectors does not. We will often denote a rule as  $(X, \bar{Y})$  where  $\bar{Y} = \{p_i\}_{i=1}^m$  if  $Y = \{p_1, \dots, p_m\}$  and  $p_1 \triangleleft \dots \triangleleft p_m$ .

A *covering*  $\mathcal{C}$  of  $\mathcal{G}$  is a sequence of rules that *covers*  $\mathcal{G}$  (i.e. the union of all the edges in  $\mathcal{C}$  equals  $E$ ). Given an edge  $e \in E$ , the *index*  $\text{index}(e)$  of  $e$  is defined to be the index of the *last* rule in the sequence  $\mathcal{C}$  that contains  $e$ . We say that  $\mathcal{C}$  is *valid* if, for all two edges  $e = (s, p), e' = (s', p')$  in  $E$  with  $e \prec e'$ , either of the following holds:

<sup>3</sup>Strictly speaking, this is only true if we are only dealing with namespace `html` (which is the case almost always in web programming and so is a reasonable assumption unless the user specifies otherwise). A node could have multiple IDs, each with a different namespace. See Section 5.

- $\text{index}(e) < \text{index}(e')$
- $\text{index}(e) = \text{index}(e')$  and, if  $(X, \{p_i\}_{i=1}^m)$  is the rule at position  $\text{index}(e)$  in  $\mathcal{C}$ , it is the case that  $p = p_j$  and  $p' = p_k$  with  $j \leq k$ .

In the example of Figure 4 with the selector `.vegetable` in the third rule, we can verify that it is indeed a valid covering of itself by observing the only ordering is  $(\text{.fruit}, \text{color:red}) \prec (\text{.vegetable}, \text{color:blue})$  and we have

$$\text{index}(\text{.fruit}, \text{color:red}) = 2 \quad \text{and} \quad \text{index}(\text{.vegetable}, \text{color:blue}) = 3.$$

This *last-occurrence* semantics reflects the cascading aspect of CSS. To relate this to the world of CSS, the original CSS file  $F$  may be represented by a CSS-graph  $\mathcal{G}$ , but  $F$  also turns out to be a valid covering of  $\mathcal{G}$ . In fact, the set of valid coverings of  $\mathcal{G}$  correspond to CSS files that are equivalent (up to reordering of selectors and property declarations) to the original CSS file. That is, if two files cover the same graph, then they will specify the same property declarations on any node of any given DOM.

To define the optimisation problem, we need to define the weight of a rule and the weight of a covering. To this end, we extend the function  $\text{wt}$  to rules and coverings by summing the weights of the nodes. More precisely, given a rule  $\bar{B} = (X, \bar{Y})$ , define

$$\text{wt}(\bar{B}) = \sum_{w \in X \cup \bar{Y}} \text{wt}(w).$$

Similarly, given a covering  $\mathcal{C} = \{\bar{B}_i\}_{i=1}^m$ , the weight  $\text{wt}(\mathcal{C})$  of  $\mathcal{C}$  is  $\sum_{i=1}^m \text{wt}(\bar{B}_i)$ . It is easy to verify that the weight of a rule (resp. covering) corresponds to the number of non-whitespace characters in a CSS rule (resp. file). The *minification problem* is, given a CSS-graph  $\mathcal{G}$ , to compute a valid covering with the minimum weight.

**(Optimal) Rule-Merging Problem** Given a CSS-graph  $\mathcal{G}$  and a covering  $\mathcal{C}$ , we define the *trim*  $\mathcal{C}_\downarrow$  of  $\mathcal{C}$  to be the covering  $\mathcal{C}'$  obtained by removing from each rule  $\bar{B} = (X, \bar{Y})$  (say at position  $i$ ) in  $\mathcal{C}$  all nodes  $v \in X \cup \bar{Y}$  that are not incident to at least one edge  $e$  with  $\text{index}(e) = i$  (i.e. the last occurrence of  $e$  in  $\mathcal{C}$ ). Such nodes  $v$  may be removed since they do not affect the validity of the covering  $\mathcal{C}$ .

We can now explain formally the file size reduction shown in Figure 5. First observe in Figure 5b that the second rule

```
.fruit, #broccoli { color:red; font-size:large }
```

has been removed. Consider the node `.fruit` and its incident edges  $(\text{.fruit}, \text{color:red})$  and  $(\text{.fruit}, \text{font-size:large})$ . Both of these edges have index 5 ( $\neq 2$ ) since they also appear in the last rule of Figure 5a. Thus we can remove the `.fruit` node from this rule. A similar argument can be made for all nodes in this rule, which means that we remove them, leaving an empty rule (not shown). In the fourth rule, the node `color:red` is incident only to  $(\text{\#tomato}, \text{color:red})$  which has index 5 ( $\neq 4$ ). The situation is the same for the node `font-size:large`, thus both of these nodes are removed from the rule.

The trim  $\mathcal{C}_\downarrow$  can be computed from  $\mathcal{C}$  and  $\mathcal{G}$  in polynomial time in a straightforward way. More precisely, first build a hashmap for the index function. Second, go through all rules  $\bar{B}$  in the covering  $\mathcal{C}$ , each node  $v$  in  $\bar{B}$ , and each edge  $e$  incident with  $v$  checking if at least one such  $e$  satisfies  $\text{index}(e) = i$ , where  $\bar{B}$  is the  $i$ th rule in  $\mathcal{C}$ . Note that  $\mathcal{C}_\downarrow$  is uniquely defined given  $\mathcal{C}$ .

We define a (*rule*)-*merging opportunity* to be a pair  $(\bar{B}, j)$  of rule  $\bar{B}$  and a number  $j \in (0, |\mathcal{C}|)$  such that  $\mathcal{C}[\bar{B} \rightarrow j]$  is a valid covering of  $\mathcal{G}$ . The *result* of applying this merging opportunity is the covering  $\mathcal{C}[\bar{B} \rightarrow j]_\downarrow$  obtained by trimming  $\mathcal{C}[\bar{B} \rightarrow j]$ . The (*rule*)-*merging problem* can be defined as follows: given a CSS-graph  $\mathcal{G}$  and a valid covering  $\mathcal{C}$ , find a merging opportunity that results in a covering with the minimum weight. This rule-merging problem is NP-hard even in the non-weighted version [72, 52].

## 5 CSS Selector Formalisation and its Intersection Problem

In this section, we will show how to efficiently compute a CSS-graph  $\mathcal{G} = (S, P, E, \prec, \text{wt})$  from a given CSS file with the help of a fast solver of quantifier-free theory of integer linear arithmetic, e.g., Z3 [18]. The key challenge is how to extract the order dependency  $\prec$  from a CSS file, which requires an algorithm for the (*selector-)*intersection problem, i.e., to check whether two given selectors can be matched by the same element in *some* document. To address this, we provide a full formalisation of CSS3 selectors [15] and a fast algorithm for the intersection problem. Since our algorithm for the intersection problem is technically very involved, we provide a short intuitive explanation behind the algorithm in this section and leave the details to Section 6.

### 5.1 Definition of Document Trees

We define the semantics of CSS3 in terms of Document Object Models (DOMs), which we also refer to as document trees. The reader may find it helpful to recall the definition of trees from Section 3.2.

A document tree consists of a number of elements, which in turn may have sub-elements as children. Each node has a *type* consisting of an element name (a.k.a. tag name) and a namespace. For example, an element `p` in the default `html` namespace is a paragraph. Namespaces commonly feature in programming languages (e.g. C++) to allow the use of multiple libraries whilst minimising the risk of overloading names. For example, the HTML designers introduced a `div` element to represent a section of an HTML document. Independent designers of an XML representation of mathematical formulas may represent division using elements also with the name `div`. Confusion is avoided by the designers specifying a namespace in addition to the required element names. In this case, the HTML designers would introduce the `div` element to the `html` namespace, while the mathematical `div` may belong to a namespace `math`. Hence, there is no confusion between `html:div` and `math:div`. As an aside, note that an HTML file may contain multiple namespaces, e.g., see [30].

Moreover, nodes may also be labelled by attributes, which take string values. For example, an HTML `img` element has a `src` attribute specifying the source of the image. Finally, a node may be labelled by a number of *pseudo-classes*. For example `:enabled` means that the node is enabled and the user may interact with it. The set of pseudo-classes is fixed by the CSS specification.

We first the formal definition before returning to the example in Section 2.1.

#### 5.1.1 Formal Definition

In the formal definition below we permit a possibly infinite set of element, namespace, and attribute names. CSS stylesheets are *not* limited to HTML documents (which have a fixed set of element names, but not attribute names since you can create custom `data-*` attributes), but they can also be applied to documents of other types (e.g. XML) that permit custom element names, namespaces, and attribute names. Thus, the sets of possible names *cannot* be fixed to finite sets from the point of view of CSS selectors, which may be applied to any document.

When it is known that the set of elements or attribute names is fixed (e.g. when only considering HTML), it is possible to adapt our approach. In particular, the small model property in Proposition 6.4 may be avoided, slightly simplifying the technique.

We denote the set of *pseudo-classes* as

$$P = \left\{ \begin{array}{l} \text{:link, :visited, :hover, :active, :focus, :target,} \\ \text{:enabled, :disabled, :checked, :root, :empty} \end{array} \right\}.$$

Then, given a possibly infinite set of *namespaces* NS, a possibly infinite set of *element names* ELE, a

possibly infinite set of *attribute names*  $A$ , and a finite alphabet<sup>4</sup>  $\Gamma$  containing the special characters  $\_$  and  $-$  (space and dash respectively), a *document tree* is a  $\Sigma$ -labelled tree  $(D, \lambda)$ , where

$$\Sigma := (\text{NS} \times \text{ELE} \times \mathcal{F}_{\text{fin}}(\text{NS} \times A, \Gamma^*) \times 2^P) .$$

Here the notation  $\mathcal{F}_{\text{fin}}(\text{NS} \times A, \Gamma^*)$  denotes the set of partial functions from  $(\text{NS} \times A)$  to  $\Gamma^*$  whose domain is finite. In other words, each node in a document tree is labeled by a namespace, an element, a function associating a finite number of namespace-attribute pairs with attribute values (strings), and zero or more of the pseudo-classes. For a function  $f_A \in \mathcal{F}_{\text{fin}}(\text{NS} \times A, \Gamma^*)$  we say  $f_A(s, a) = \perp$  when  $f_A$  is undefined over  $s \in \text{NS}$  and  $a \in A$ , where  $\perp \notin \Gamma^*$  is a special undefined value. Furthermore, we assume special attribute names `class`, `id`  $\in A$  that will be used to attach classes (see later) and IDs to nodes.

When  $\lambda(\eta) = (s, e, f_A, P)$  we define the following projections of the labelling function

$$\begin{aligned} \lambda_S(\eta) &= s, \\ \lambda_E(\eta) &= e, \\ \lambda_A(\eta) &= f_A, \text{ and} \\ \lambda_P(\eta) &= P. \end{aligned}$$

We will use the following standard XML notation: for an element name  $e$ , a namespace  $s$ , and an attribute name  $a$ , let  $s:e$  (resp.  $s:a$ ) denote the pair  $(s, e)$  (resp.  $(s, a)$ ). The notation helps to clarify the role of namespaces as a means of providing contextual or scoping information to an element/attribute.

There are several consistency constraints on the node labellings.

- For each  $s \in \text{NS}$ , there are *no* two nodes in the tree with the same value of  $s:\text{id}$ .
- A node cannot be labelled by both `:link` and `:visited`.
- A node cannot be labelled by both `:enabled` and `:disabled`.
- Only one node in the tree may be labelled `:target`.
- A node contains the label `:root` iff it is the root node.
- A node labelled `:empty` must have no children.

From now on, we will tacitly assume that document trees satisfy these consistency constraints. We write  $\text{Trees}(\text{NS}, \text{ELE}, A, \Gamma)$  for the set of such trees.

### 5.1.2 Example

Consider the HTML file in Figure 1. This file can be represented precisely by the tree in Figure 7. In this tree, each node is first labelled by its type, which consists of its namespace and its element name. In all cases, the namespace is the `html` namespace, while the element names are exactly as in the HTML file. In addition, we label each node with the attributes and pseudo-classes with which they are labelled. The `html:html` node is labelled with `:root` since it is the root element of the tree. The `html:div` node is labelled with the ID `heading`. The `html:img` node is labelled with the attribute `class` with value `banner` and the attribute `html:src` with value `banner.jpg`, as well as the pseudo-class `:empty` indicating that the node has no contents. The remaining nodes however are not labelled `:empty` since even the leaf nodes contain some text (which is not represented in our tree model as it is not matchable by a selector). Hence, a node with no child may still be non-empty.

<sup>4</sup>See the notes at the end of the section.

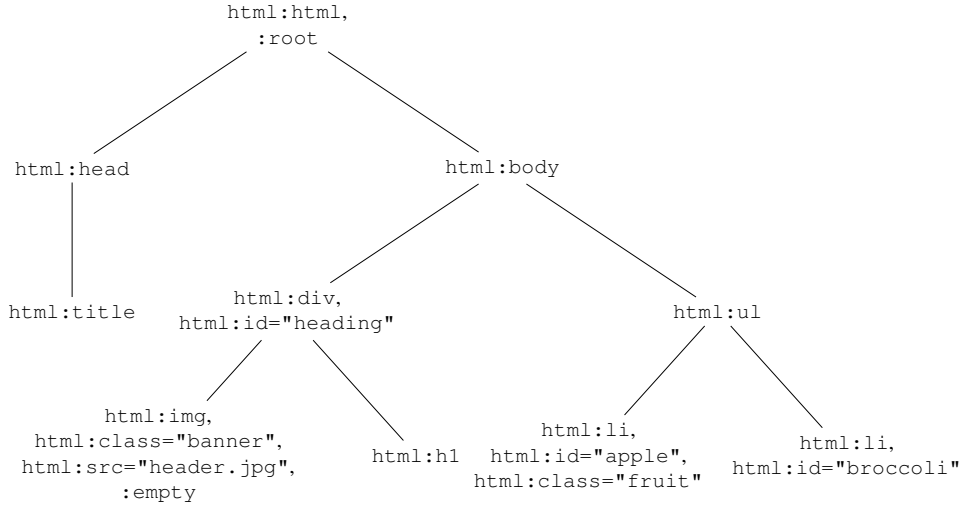


Figure 7: A DOM tree representation of the HTML file.

Formally, the DOM tree is  $(D, \lambda)$  where  $D = \{\varepsilon, 1, 2, 11, 21, 22, 211, 212, 221, 222\}$  and

$$\begin{aligned}
 \lambda(\varepsilon) &= (\text{html}, \text{html}, \emptyset, \{ : \text{root} \}) \\
 \lambda(1) &= (\text{html}, \text{head}, \emptyset, \emptyset) \\
 \lambda(2) &= (\text{html}, \text{body}, \emptyset, \emptyset) \\
 \lambda(11) &= (\text{html}, \text{title}, \emptyset, \emptyset) \\
 \lambda(21) &= (\text{html}, \text{div}, (\text{html}, \text{id}) \mapsto \text{heading}, \emptyset) \\
 \lambda(22) &= (\text{html}, \text{ul}, \emptyset, \emptyset) \\
 \lambda(211) &= \left( \text{html}, \text{img}, \left( \begin{array}{l} (\text{html}, \text{class}) \mapsto \text{banner}, \\ (\text{html}, \text{src}) \mapsto \text{header.jpg} \end{array} \right), \{ : \text{empty} \} \right) \\
 \lambda(212) &= (\text{html}, \text{h1}, \emptyset, \emptyset) \\
 \lambda(221) &= \left( \text{html}, \text{li}, \left( \begin{array}{l} (\text{html}, \text{id}) \mapsto \text{apple}, \\ (\text{html}, \text{class}) \mapsto \text{fruit} \end{array} \right), \emptyset \right) \\
 \lambda(222) &= (\text{html}, \text{li}, (\text{html}, \text{id}) \mapsto \text{broccoli}, \emptyset) .
 \end{aligned}$$

## 5.2 Definition of CSS3 selectors

In the following sections we define CSS selectors syntax and semantics. Informally, a CSS selector consists of *node selectors*  $\sigma$  — which match individual nodes in the tree — combined using the operators  $\gg$ ,  $>$ ,  $+$ , and  $\sim$ . These operators express the descendant-of, child-of, neighbour-of, and sibling-of relations respectively. Note that the blank space character is used instead of  $\gg$  in CSS3, though we opt for the latter in the formalisation for the sake of readability. So, for example, we use `.journal >> .science` (i.e. choose all nodes with class `.science` that is a descendant of nodes with class `.journal`) instead of the standard syntax `.journal .science`. In addition, in order to distinguish syntax from meaning, we use slightly different notation to their counterpart semantical operators  $\otimes$ ,  $\odot$ ,  $\oplus$ , and  $\odot$ .

We remark that a comma  $(,)$  is not an operator in the CSS selector syntax. Instead a *selector group* is a comma-separated list of selectors that is matched if any of its selectors is matched. A CSS rule thus consists of a selector group and a list of property declarations. For the purposes of rule-merging it is desirable to treat selectors individually as it allows the most flexibility in reorganising the CSS file. Hence we treat a selector group simply as a set of selectors that can be separated if needed.

A node selector  $\sigma$  has the form  $\tau\Theta$  where  $\tau$  constrains the *type* of the node. That is,  $\tau$  places



restrictions on the element label of the node, e.g., `p` for paragraph elements and `*` (or an empty string) for all elements. The rest of the selector is a set  $\Theta$  of *simple* selectors (written as a concatenation of strings representing these simple selectors) that assert atomic properties of the node. There are four types of simple selectors.

*Type 1:* attribute selectors of the form  $[s|a\ op\ v]$  for some namespace  $s$ , attribute  $a$ , operator  $op \in \{=, \sim, |=, ^=, \$=, *=\}$ , and some string  $v \in \Gamma^*$ . We may write  $[a\ op\ v]$  to mean that a node can be matched by  $[s|a\ op\ v]$  for some  $s$ . The operators  $=$ ,  $^=$ ,  $\$=$ , and  $*=$  take their meaning from regular expressions. That is, equals, begins-with, ends-with, and contains respectively. The remaining operators are more subtle. The  $\sim$  operator means the attribute is a string of space-separated values, one of which is  $v$ . The  $|=$  operator is intended for use with language identification, e.g., as in the attribute selector  $[\text{lang}\ |=\ \text{"en-GB"}]$  to mean “English” as spoken in Great Britain. Thus  $|=$  asserts that either the attribute has value  $v$  or is a string of the form  $v-v'$  where  $-$  is the dash character, and  $v'$  is some string. Note that if the `lang` attribute value of a node is `en-GB`, the node also matches the simple selector  $[\text{lang}\ |=\ \text{"en"}]$ . In addition, recall that `class` and `id` are two special attribute names. For this reason, CSS introduces the shorthands  $.v$  and  $\#v$  for, respectively, the simple selectors  $[\text{class}\ \sim\ v]$  and  $[\text{id}\ =\ v]$ , i.e., asserting that the node has a class or ID  $v$ . An example of a valid CSS selector is the selector `h1.fruit.vegetable`, which chooses all nodes with class `fruit` and `vegetable`, and element name `h1` (which includes the following two elements: `<h1 class="fruit vegetable">` and `<h1 class="vegetable fruit">`).

*Type 2:* attribute selectors of the form  $[s|a]$ , asserting that the attribute is merely defined on the node. As before, we may write  $[a]$  to mean that the node may be matched by  $[s|a]$  for some namespace  $s$ . As an example, `img[alt]` chooses all `img` elements where the attribute `alt` is defined.

*Type 3:* pseudo-class label of a node, e.g., the selector `:enabled` ensures the node is currently enabled in the document. There are several further kinds of pseudo-classes that assert counting constraints on the children of a selected node. Of particular interest are selectors such as `:nth-child( $\alpha n + \beta$ )`, which assert that the node has a particular position in the sibling order. For example, the selector `:nth-child( $2n + 1$ )` means there is some  $n \geq 0$  such that the node is the  $(2n + 1)$ st node in the sibling order.

*Type 4:* negations `:not( $\theta$ )` of a simple selector  $\theta$  with the condition that negations cannot be nested or apply to multiple atoms. For example, `:not(.fruit):not(.vegetable)` is a valid selector, whereas `:not(:not(.vegetable))` and `:not(.fruit.vegetable)` are *not* a valid selectors.

### 5.2.1 Syntax

Fix the sets  $\text{NS}$ ,  $\text{ELE}$ ,  $A$ , and  $\Gamma$ . We define  $\text{SEL}$  for the set of (CSS) *selectors* and  $\text{NSEL}$  for the set of *node selectors*. In the following  $|$  will be used to separate syntax alternatives, while  $\cdot$  is an item of CSS syntax. The set  $\text{SEL}$  is the set of formulas  $\varphi$  defined as:

$$\varphi ::= \sigma \mid \varphi \gg \sigma \mid \varphi > \sigma \mid \varphi + \sigma \mid \varphi \sim \sigma$$

where  $\sigma \in \text{NSEL}$  is a *node selector* with syntax  $\sigma ::= \tau \Theta$  with  $\tau$  having the form

$$\tau ::= * \mid (s|*) \mid e \mid (s|e)$$

where  $s \in \text{NS}$  and  $e \in \text{ELE}$  and  $\Theta$  is a (possibly empty) set of conditions  $\theta$  with syntax

$$\theta ::= \theta_{\neg} \mid \text{:not}(\sigma_{\neg})$$

where  $\theta_{\neg}$  and  $\sigma_{\neg}$  are conditions that do not contain negations, i.e.:

$$\begin{aligned}
\sigma_{\neg} &::= * \mid (s \mid *) \mid e \mid (s \mid e) \mid \theta_{\neg} \\
\theta_{\neg} &::= [s \mid a] \mid [s \mid a \text{ op } v] \mid [a] \mid [a \text{ op } v] \mid \\
&\quad :link \mid :visited \mid :hover \mid :active \mid :focus \mid \\
&\quad :enabled \mid :disabled \mid :checked \mid \\
&\quad :root \mid :empty \mid :target \mid \\
&\quad :nth\text{-child}(\alpha n + \beta) \mid :nth\text{-last-child}(\alpha n + \beta) \mid \\
&\quad :nth\text{-of-type}(\alpha n + \beta) \mid :nth\text{-last-of-type}(\alpha n + \beta) \\
&\quad :only\text{-child} \mid :only\text{-of-type} \\
op &::= = \mid \sim = \mid |= \mid ^ = \mid \$ = \mid * =
\end{aligned}$$

where  $s \in \text{NS}$ ,  $e \in \text{ELE}$ ,  $a \in A$ ,  $v \in \Gamma^*$ , and  $\alpha, \beta \in \mathbb{Z}$ . Whenever  $\Theta$  is the empty set, we will denote the node selector  $\tau\Theta$  as  $\tau$  instead of  $\tau\emptyset$ .

### 5.2.2 Semantics

The semantics of a selector is defined with respect to a document tree and a node in the tree. More precisely, the semantics of CSS3 selectors  $\varphi$  are defined inductively with respect to a document tree  $T = (D, \lambda)$  and a node  $\eta \in D$  as follows. (Note: (1)  $p$  ranges over the set  $P$  of pseudo-classes, (2)  $vv'$  is the concatenation of the strings  $v$  and  $v'$ , and (3)  $v\text{-}v'$  is the concatenation of  $v$  and  $v'$  with a “-” in between.)

$$\begin{aligned}
T, \eta \models \varphi \gg \sigma &\stackrel{\text{def}}{\Leftrightarrow} \exists \eta' \otimes \eta . (T, \eta' \models \varphi) \text{ and } (T, \eta \models \sigma) \\
T, \eta \models \varphi > \sigma &\stackrel{\text{def}}{\Leftrightarrow} \exists \eta' \odot \eta . (T, \eta' \models \varphi) \text{ and } (T, \eta \models \sigma) \\
T, \eta \models \varphi + \sigma &\stackrel{\text{def}}{\Leftrightarrow} \exists \eta' \oplus \eta . (T, \eta' \models \varphi) \text{ and } (T, \eta \models \sigma) \\
T, \eta \models \varphi \sim \sigma &\stackrel{\text{def}}{\Leftrightarrow} \exists \eta' \ominus \eta . (T, \eta' \models \varphi) \text{ and } (T, \eta \models \sigma) \\
T, \eta \models \tau\Theta &\stackrel{\text{def}}{\Leftrightarrow} (T, \eta \models \tau) \text{ and } \forall \theta \in \Theta . (T, \eta \models \theta) \\
T, \eta \models (s \mid *) &\stackrel{\text{def}}{\Leftrightarrow} s = \lambda_S(\eta) \\
T, \eta \models * &\stackrel{\text{def}}{\Leftrightarrow} \top \\
T, \eta \models (s \mid e) &\stackrel{\text{def}}{\Leftrightarrow} s = \lambda_S(\eta) \wedge e = \lambda_E(\eta) \\
T, \eta \models e &\stackrel{\text{def}}{\Leftrightarrow} T, \eta \models (s \mid e) \text{ for some } s \in \text{NS} \\
T, \eta \models p &\stackrel{\text{def}}{\Leftrightarrow} p \in \lambda_P(\eta) \\
T, \eta \models \text{:not } (\theta_{\neg}) &\stackrel{\text{def}}{\Leftrightarrow} \neg (T, \eta \models \theta_{\neg}) \\
T, \eta \models [a] &\stackrel{\text{def}}{\Leftrightarrow} T, \eta \models [s \mid a] \text{ for some } s \in \text{NS} \\
T, \eta \models [a \text{ op } v] &\stackrel{\text{def}}{\Leftrightarrow} T, \eta \models [s \mid a \text{ op } v] \text{ for some } s \in \text{NS} \\
T, \eta \models [s \mid a] &\stackrel{\text{def}}{\Leftrightarrow} \lambda_A(\eta)(s, a) \neq \perp \\
T, \eta \models [s \mid a = v] &\stackrel{\text{def}}{\Leftrightarrow} \lambda_A(\eta)(s, a) = v \\
T, \eta \models [s \mid a \text{ |= } v] &\stackrel{\text{def}}{\Leftrightarrow} \left( \lambda_A(\eta)(s, a) = v \text{ or } \exists v' . (\lambda_A(\eta)(s, a) = v\text{-}v') \right)
\end{aligned}$$

$$\begin{aligned}
T, \eta \models [s | a \wedge = v] &\stackrel{\text{def}}{\Leftrightarrow} \exists v' \in \Gamma^* . \lambda_A(\eta)(s, a) = vv' \\
T, \eta \models [s | a \$ = v] &\stackrel{\text{def}}{\Leftrightarrow} \exists v' \in \Gamma^* . \lambda_A(\eta)(s, a) = v'v \\
T, \eta \models [s | a * = v] &\stackrel{\text{def}}{\Leftrightarrow} \exists v_1, v_2 \in \Gamma^* . \lambda_A(\eta)(s, a) = v_1vv_2
\end{aligned}$$

with the missing attribute selector being (noting  $v \_ v'$  is the concatenation of  $v$  and  $v'$  with the space character  $\_$  in between)

$$\begin{aligned}
T, \eta \models [a \sim = v] &\stackrel{\text{def}}{\Leftrightarrow} \lambda_A(\eta)(s, a) = v \text{ or } \exists v' . (\lambda_A(\eta)(s, a) = v \_ v') \text{ or} \\
&\quad \exists v' . (\lambda_A(\eta)(s, a) = v' \_ v) \text{ or } \exists v_1, v_2 . (\lambda_A(\eta)(s, a) = v_1 \_ v \_ v_2)
\end{aligned}$$

then, for the counting selectors

$$\begin{aligned}
T, \eta \models \text{:nth-child}(\alpha n + \beta) &\stackrel{\text{def}}{\Leftrightarrow} \text{there is some } n \in \mathbb{N} \text{ such that } \eta \text{ is the } \alpha n + \beta \text{th} \\
&\quad \text{child} \\
T, \eta \models \text{:only-child} &\stackrel{\text{def}}{\Leftrightarrow} \text{the parent of } \eta \text{ has precisely one child} \\
T, \eta \models \text{:nth-of-type}(\alpha n + \beta) &\stackrel{\text{def}}{\Leftrightarrow} \text{there is some } n \in \mathbb{N} \text{ such that the parent of } \eta \text{ has} \\
&\quad \text{precisely } \alpha n + \beta - 1 \text{ children with namespace } \lambda_S(\eta) \\
&\quad \text{and element name } \lambda_E(\eta) \text{ for some } n \text{ that are (strictly)} \\
&\quad \text{preceding siblings of } \eta \\
T, \eta \models \text{:only-of-type} &\stackrel{\text{def}}{\Leftrightarrow} \text{the parent of } \eta \text{ has precisely one child with} \\
&\quad \text{namespace } \lambda_S(\eta) \text{ and element name } \lambda_E(\eta)
\end{aligned}$$

Finally, the semantics of the remaining two selectors, which are  $\text{:nth-last-child}(\alpha n + \beta)$  and  $\text{:nth-last-of-type}(\alpha n + \beta)$ , is exactly the same as  $\text{:nth-child}(\alpha n + \beta)$  and  $\text{:nth-of-type}(\alpha n + \beta)$ , respectively, except with the sibling ordering reversed (i.e. the rightmost child of a parent is treated as the first).

**Remark 5.1.** *Readers familiar with HTML may have expected more constraints in the semantics. For example, if a node matches  $\text{:hover}$ , then its parent should also match  $\text{:hover}$ . However, this is part of the HTML5 specification, not of CSS3. In fact, the CSS3 selectors specification explicitly states that a node matching  $\text{:hover}$  does not imply its parent must also match  $\text{:hover}$ .*

### 5.2.3 Divergences from full CSS

Note that we diverge from the full CSS specification in a number of places. However, we do not lose expressivity.

- We assume each element has a namespace. In particular, we do not allow elements without a namespace. There is no loss of generality here since we can simply assume a “null” namespace is used instead. Moreover, we do not support default name spaces and assume namespaces are explicitly given.
- We did not include  $\text{:lang}(l)$ . Instead, we will assume (for convenience) that all nodes are labelled with a language attribute with some fixed namespace  $s$ . In this case,  $\text{:lang}(l)$  is equivalent<sup>5</sup> to  $[s | \text{lang} = l]$ .

<sup>5</sup>The CSS specification defines  $\text{:lang}(l)$  in this way. A restriction of the language values to standardised language codes is only a recommendation.

- We did not include `:indeterminate` since it is not formally part of the CSS3 specification.
- We omit the selectors `:first-child` and `:last-child`, as well as `:first-of-type` and `:last-of-type`, since they are expressible using the other operators.
- We omitted `even` and `odd` from the `nth` child operators since these are easily definable as  $2n$  and  $2n + 1$ .
- We do not explicitly handle document fragments. These may be handled in a number of ways. For example, by adding a phantom root element (since the root of a document fragment does not match `:root`) with a fresh ID  $\iota$  and adjusting each node selector in the CSS selector to assert `:not(# $\iota$ )`. Similarly, lists of document fragments can be modelled by adding several subtrees to the phantom root.
- A CSS selector can be suffixed with a *pseudo-element* of the form `::first-line`, `::first-letter`, `::before`, and `::after`. Pseudo-elements are easy to handle and only provide a distraction to our presentation. For this reason, we relegate them into the appendix.
- We define our DOM trees to use a finite alphabet  $\Gamma$ . Currently the CSS3 selectors specification uses Unicode as its alphabet for lexing. Although the CSS3 specification is not explicit about the finiteness of characters appearing in potential DOMs, since Unicode is finite [70] (with a maximal possible codepoint) we feel it is reasonable to assume DOMs are also defined over a finite alphabet.

### 5.3 Solving the intersection problem

We now address the problem of checking the intersection of two CSS selectors. Let us write

$$\llbracket \varphi \rrbracket := \{(T, \eta) : T, \eta \models \varphi\}$$

to denote the set of pairs of tree and node satisfying the selector  $\varphi$ . The *intersection problem of CSS selectors* is to decide if  $\llbracket \varphi \rrbracket \cap \llbracket \varphi' \rrbracket \neq \emptyset$ , for two given selectors  $\varphi$  and  $\varphi'$ . A closely related decision problem is the *non-emptiness problem of CSS selectors*, which is to decide if  $\llbracket \varphi \rrbracket \neq \emptyset$ , for a given selector  $\varphi$ . The two problems are *not* the same since CSS selectors are not closed under intersection (i.e. the conjunction of two CSS selectors is in general not a valid CSS selector).

**Theorem 5.2** (Non-Emptiness). *The non-emptiness problem for CSS selectors is efficiently reducible to satisfiability over quantifier-free theory over integer linear arithmetic. Moreover, the problem is NP-complete.*

**Theorem 5.3** (Intersection). *The intersection problem for CSS selectors is efficiently reducible to satisfiability over quantifier-free theory over integer linear arithmetic. Moreover, the problem is NP-complete.*

Recall from Section 3 that satisfiability over quantifier-free theory over integer linear arithmetic is in NP and can be solved by a highly-optimised SMT solver (e.g. Z3 [18]). The NP-hardness in the above theorems suggests that our SMT-based approach is theoretically optimal. In addition, our experiments with real-world selectors (see Section 8) suggest that our SMT-based approach can be optimised to be fast in practice (see Appendix D.1), with each problem instance solved within moments.

**Idea behind our SMT-based approach** We now provide the idea behind the efficient reduction to quantifier-free theory over integer linear arithmetic. Our reduction first goes via a new class of tree automata (called CSS automata), which — like CSS selectors — are symbolic representations of sets of pairs containing a document tree and a node in the tree. We will call such sets *languages recognised by the automata*. Given a CSS selector  $\varphi$ , we can efficiently construct a CSS automaton  $\mathcal{A}$  that can symbolically represent  $\llbracket \varphi \rrbracket$ . Unlike CSS selectors, however, we will see that languages recognised by CSS automata enjoy closure under intersection, which will allow us to treat the intersection problem as the non-emptiness problem. More precisely, a CSS automaton navigates a tree structure in a similar manner to a CSS selector: transitions may only move down the tree or to a sibling, while checking a number of properties on the visited nodes. The difficulty of taking a direct intersection of two selectors is that the two selectors may descend to different child nodes, and then meet again after the application of a number of sibling combinators, i.e., their paths may diverge and combine several times. CSS automata overcome this difficulty by always descending to the *first* child, and then move from *sibling to sibling*. Thus, the intersection of CSS automata can be done with a straightforward automata product construction, e.g., see [71].

Next, we show that the non-emptiness of CSS automata can be decided in NP by a polynomial-time reduction to satisfiability of quantifier-free theory of integer linear arithmetic. In our experiments, we used an optimised version of the reduction, which is detailed in Appendix D.1. For non-emptiness, ideally, we would like to show that if a CSS automaton has a non-empty language, then it accepts a small tree (i.e. with polynomially many nodes). This is unfortunately not the case, as the reader can see in our NP-hardness proof idea below. Therefore, we use a different strategy. First, we prove three “small model lemmas”. The first is quite straightforward and shows that, to prove non-emptiness, it suffices to consider a witnessing automata run of length  $n$  for an automaton with  $n$  transitions (each automata transition allows some nodes to be skipped). Second, we show that it suffices to consider attribute selector values (i.e. strings) of length linear in the size of the CSS automata. This is non-trivial and uses a construction inspired by [45]. Third, we show that it suffices to consider trees whose sets of namespaces and element names are linear in the size of the CSS automaton. Our formula  $\varphi$  attempts to guess this automata run, the attribute selector values, element names, and namespaces. The global ID constraint (i.e. all the guessed IDs are distinct) can be easily asserted in the formula. So far, boolean variables are sufficient because the small model lemmas allow us to do bit-blasting. *Where, then, do the integer variables come into play?* For each position  $i$  in the guessed path, we introduce an integer variable  $\bar{n}_i$  to denote that the node at position  $i$  in the path is the  $\bar{n}_i$ th child. This is necessary if we want to assert counting constraints like  $\text{:nth-child}(\alpha n + \beta)$ , which would be encoded in integer linear arithmetic as  $\exists \bar{n} : \bar{n}_i = \alpha \bar{n} + \beta$ .

**Proof Idea of NP-hardness** We now provide an intuition on how to prove NP-hardness in the above theorems. First, observe that the intersection is computationally at least as hard as the non-emptiness problem since we can set the second selector to be  $*$ . To prove NP-hardness of the non-emptiness, we give a polynomial-time reduction from the NP-complete problem of *non-universality of unions of arithmetic progressions* [65, Proof of Theorem 6.1]. Examples of arithmetic progressions are  $2\mathbb{N} + 1$  and  $5\mathbb{N} + 2$ , which are shorthands for the sets  $\{1, 3, 5, \dots\}$  and  $\{2, 7, 12, \dots\}$ , respectively. The aforementioned non-universality problem allows an arbitrary number of arithmetic progressions as part of the input and we are interested in checking whether the union equals the entire set  $\mathbb{N}$  of natural numbers. As an example of the reduction, checking  $\mathbb{N} \neq 2\mathbb{N} + 1 \cup 5\mathbb{N} + 2$  is equivalent to the non-emptiness of

$$\text{:not}(\text{root}) : \text{not}(\text{:nth-child}(2n+2)) : \text{not}(\text{:nth-child}(5n+3))$$

which can be paraphrased as checking the existence of a tree with a node that is neither the root, nor the

$2n + 2$ nd child, nor the  $5n + 3$ rd child. Observe that we add 1 to the offset of the arithmetic progressions since the selector `:nth-child` starts counting (the number of children) from 1, not from 0. A full NP-hardness proof is available in Appendix B.2.

#### 5.4 Extracting the edge order $\prec$ from a CSS file

Recall that our original goal is to compute a CSS-graph  $\mathcal{G} = (S, P, E, \prec, \text{wt})$  from a given CSS file. The sets  $P$ ,  $S$ , and  $E$ , and the function  $\text{wt}$  can be computed easily as explained in Section 4. We now show how to compute  $\prec$  using the algorithm for checking intersection of two selectors. We present an intuitive ordering, before explaining how this may be relaxed while still preserving the semantics.

An initial definition of  $\prec$  is simple to obtain: we want to order  $(s, p) \prec (s', p')$  whenever  $(s', p')$  appears later in the CSS file than  $(s, p)$ , the selectors may overlap but are not distinguished by their specificity, and  $p$  and  $p'$  assign conflicting values to a given property name. More formally, we first compute the specificity of all the selectors in the CSS file. This can be easily computed in the standard way [15]. Now, the relation  $\prec$  can only relate two edges  $(s, p), (s', p') \in E$  satisfying

1.  $s$  and  $s'$  have the same specificity,
2. we have  $p \neq p'$  but the property names for  $p$  and  $p'$  are the same (e.g.  $p = \text{color:blue}$  and  $p' = \text{color:red}$  with property name `color`), and
3.  $s$  intersects with  $s'$  (i.e.  $\llbracket s \rrbracket \cap \llbracket s' \rrbracket \neq \emptyset$ ).

If both (1) and (2) are satisfied, Condition (3) can be checked by means of SMT-solver via the reduction in Theorem 5.3. Supposing that Condition (3) holds, we simply compute the indices of the edges in the file:  $m := \text{index}((s, p))$  and  $m' := \text{index}((s', p'))$ . Recall  $\text{index}(e)$  was defined formally in Section 4. We put  $(s, p) \prec (s', p')$  iff  $m < m'$ . There are two minor technical details with the keyword `!important` and *shorthand property names*; see Appendix B.3.

The ordering given above unfortunately turns out to be too conservative. In the following, we give an example to demonstrate this, and propose a refinement to the ordering. Consider the CSS file

```
.a { color:red; color:rgba(255,0,0,0.5) }
.b { color:red; color:rgba(255,0,0,0.5) }
```

In this file, both nodes matching `.a` and `.b` are assigned a semi-transparent red with solid red being defined as a *fallback* when the semi-transparent red is not supported. If the edge order is calculated as above, we obtain

$$(.a, \text{color:rgba}(255, 0, 0, 0.5)) \prec (.b, \text{color:red}) \quad (1)$$

which prevents the obvious rule-merging

```
.a, .b { color:red; color:rgba(255,0,0,0.5) }
```

The key observation is that the fact that we also have

$$(.b, \text{color:red}) \prec (.b, \text{color:rgba}(255, 0, 0, 0.5)) \quad (2)$$

renders any violations of (1) benign: such a violation would give precedence to the declaration `color:rgba(255, 0, 0, 0.5)` over `color:red` for nodes matching both `.a` and `.b`. However, because of (2) this should happen anyway and we can omit (1) from  $\prec$ .

Formally, the ordering we need is as follows. If Conditions (1-3) hold, we compute  $m := \text{index}((s, p))$  and  $m' := \text{index}((s', p'))$  and put  $(s, p) \prec (s', p')$  iff

- $m < m'$ , and
- $(s', p)$  does not exist or  $\text{index}((s', p)) < m'$ .

That is, we omit  $(s, p) \prec (s', p')$  if  $(s', p)$  appears later in the CSS file (i.e.  $\text{index}((s', p')) < \text{index}((s', p))$ ). Note, we are guaranteed in this latter case to include  $(s', p') \prec (s', p)$  since  $(s', p')$  and  $(s', p)$  can easily be seen to satisfy the conditions for introducing  $(s', p') \prec (s', p)$ .

## 6 More Details on Solving Selector Intersection Problem

In the previous section, we have given the intuition behind the efficient reduction from the CSS selector intersection problem to quantifier-free theory over integer linear arithmetic, for which there is a highly-optimised SMT-solver [18]. In this section, we present this reduction in full, which may be skipped on a first reading without affecting the flow of the paper.

This section is structured as follows. We begin by defining CSS automata. We then provide a semantic-preserving transformation from CSS selectors to CSS automata. Next we show the closure of CSS automata under intersection. The closure allows us to reduce the intersection problem of CSS automata to the non-emptiness problem of CSS automata. Finally, we provide a reduction from non-emptiness of CSS automata to satisfiability over quantifier-free integer linear arithmetic. We will see that each such transformation/reduction runs in polynomial-time, resulting in the complexity upper bound of NP, which is precise due to the NP-hardness of the problem from the previous section.

### 6.1 CSS Automata

CSS automata are a kind of finite automata which navigate the tree structure of a document tree. Transitions of the automata will contain one of four labels:  $\downarrow$ ,  $\rightarrow$ ,  $\rightarrow_+$ , and  $\circ$ . Intuitively, these transitions perform the following operations.  $\downarrow$  moves to the first child of the current node.  $\rightarrow$  moves to the next sibling of the current node.  $\rightarrow_+$  moves an arbitrary number of siblings to the right. Finally,  $\circ$  reads the node matched by the automaton. Since CSS does not have loops, we require only self loops in our automata, which are used to skip over nodes (e.g.  $.v \sim .v'$  may pass over many nodes between those matching  $.v$  and those matching  $.v'$ ). We do not allow  $\rightarrow$  to label a loop – this is for the purposes of the NP proof: it can be more usefully represented as  $\rightarrow_+$ .

An astute reader may complain that  $\rightarrow_+$  does not need to appear on a loop since it can already pass over an arbitrary number of nodes. However, the product construction used for intersection becomes easier if  $\rightarrow_+$  appears only on loops. There is no analogue of  $\rightarrow_+$  for  $\downarrow$  because we do not need it: the use of  $\rightarrow_+$  is motivated by selectors such as  $:\text{nth-child}(\alpha n + \beta)$  which count the number of siblings of a node. No CSS selector counts the number of descendants/ancestors.

**Formal Definition of CSS Automata** A CSS Automaton  $\mathcal{A}$  is a tuple  $(Q, \Delta, q^{\text{in}}, q_f)$  where  $Q$  is a finite set of states,  $\Delta \subseteq Q \times \{\downarrow, \rightarrow, \rightarrow_+, \circ\} \times \text{NSEL} \times Q$  is a transition relation,  $q^{\text{in}} \in Q$  is the initial state, and  $q_f \in Q$  is the final state. Moreover,

1. (only self-loops) there exists a partial order  $\lesssim$  such that  $(q, d, \sigma, q') \in \Delta$  implies  $q' \lesssim q$ ,
2. ( $\rightarrow_+$  loops and doesn't check nodes) for all  $(q, \rightarrow_+, \sigma, q') \in \Delta$  we have  $q = q'$  and  $\sigma = *$ ,
3. ( $\rightarrow$  doesn't label loops) for all  $(q, d, \sigma, q) \in \Delta$  we have  $d \neq \rightarrow$  and  $\sigma = *$ ,
4. ( $\circ$  checks last node only) for all  $(q, d, \sigma, q') \in \Delta$  we have  $q' = q_f$  iff  $d = \circ$ , and

5. ( $q_f$  is a sink) for all  $(q, d, \sigma, q') \in \Delta$  we have  $q \neq q_f$ .

We now define the semantics of CSS automata, i.e., given an automaton  $\mathcal{A}$ , which language  $\mathcal{L}(\mathcal{A})$  they recognise. Intuitively, the set  $\mathcal{L}(\mathcal{A})$  contains the set of pairs of document tree and node, which the automaton  $\mathcal{A}$  accepts. We will now define this more formally. Write  $q \xrightarrow[\sigma]{d} q'$  to denote a transition  $(q, d, \sigma, q') \in \Delta$ . A document tree  $T = (D, \lambda)$  and node  $\eta \in D$  is *accepted* by a CSS automaton  $\mathcal{A}$  if there exists a sequence

$$q_0, \eta_0, q_1, \eta_1, \dots, q_\ell, \eta_\ell, q_{\ell+1} \in (Q \times D)^* \times \{q_f\}$$

such that  $q_0 = q^{\text{in}}$  is the initial state,  $\eta_0 = \varepsilon$  is the root node,  $q_{\ell+1} = q_f$  is the final state,  $\eta_\ell = \eta$  is the matched node, and for all  $i$ , there is some transition  $q_i \xrightarrow[\sigma]{d} q_{i+1}$  with  $\eta_i$  satisfying  $\sigma$  and if  $i \leq \ell$ ,

1. if  $d = \downarrow$  then  $\eta_{i+1} = \eta_i 1$  (i.e., the leftmost child of  $\eta_i$ ),
2. if  $d = \rightarrow$  then there is some  $\eta'$  and  $\iota$  such that  $\eta_i = \eta' \iota$  and  $\eta_{i+1} = \eta'(\iota + 1)$ , and
3. if  $d = \rightarrow_+$  then there is some  $\eta'$ ,  $\iota$  and  $\iota'$  such that  $\eta_i = \eta' \iota$  and  $\eta_{i+1} = \eta' \iota'$  and  $\iota' > \iota$ .

Such a sequence is called an *accepting run* of length  $\ell$ . The *language*  $\mathcal{L}(\mathcal{A})$  *recognised* by  $\mathcal{A}$  is the set of pairs  $(T, \eta)$  accepted by  $\mathcal{A}$ .

## 6.2 Transforming CSS Selectors to CSS Automata

The following proposition shows that CSS automata are no less expressive than CSS selectors.

**Proposition 6.1.** *For each CSS selector  $\varphi$ , we may construct in polynomial-time a CSS automaton  $\mathcal{A}_\varphi$  such that  $\mathcal{L}(\mathcal{A}_\varphi) = \llbracket \varphi \rrbracket$ .*

We show this proposition by giving a translation from a given CSS selector to a CSS automaton. Before the formal definition, we consider the a simple example. A more complex example is shown after the translation.

### 6.2.1 Simple Example

Consider the selector

$$p + .a$$

which selects a node that has a class  $a$  and is directly a right neighbour of a node with element  $p$ . Figure 8 gives a CSS automaton representing the selector. The automaton begins with a loop that can navigate down any branch of the tree using the  $\downarrow$  and  $\rightarrow_+$  transitions from  $o_1$ . Then, since it always moves from the first child to the last, it will first see the node with the  $p$ . When reading this node, it will move to the next child using  $\rightarrow$  before matching the node with class  $a$ , leading to the accepting state.

### 6.2.2 Formal Translation

Given a CSS selector  $\varphi$ , we define  $\mathcal{A}_\varphi$  as follows. We can write  $\varphi$  uniquely in the form

$$\sigma_1 o_1 \sigma_2 o_2 \cdots o_{n-1} \sigma_n$$

where each  $\sigma_i$  is a node selector, and each  $o_i \in \{\gg, >, +, \cdot\}$ . We will have a state  $o_i$  corresponding to each  $\sigma_i, o_i$ . We define

$$\mathcal{A}_\varphi = (Q, \text{ELE}, \Delta, o_1, q_f)$$



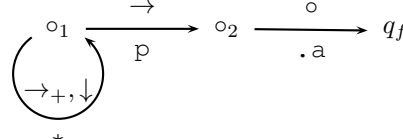


Figure 8: CSS Automaton for  $p + . a$ .

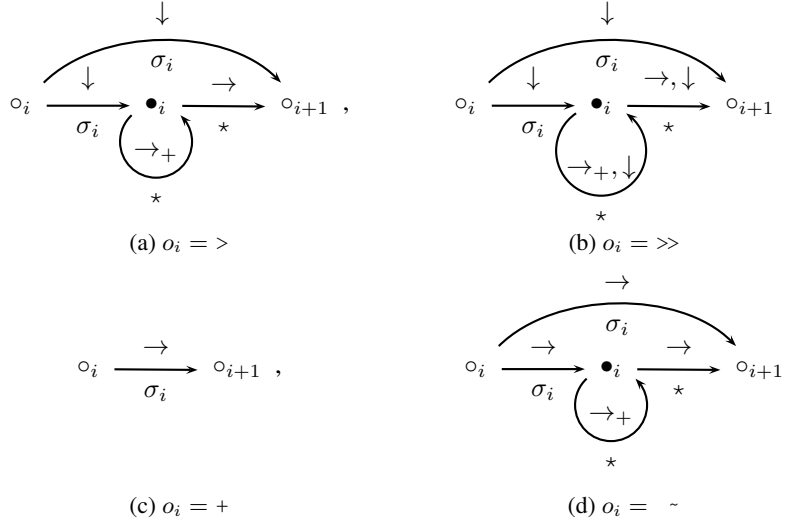


Figure 9: Converting selectors to CSS automata.

where  $Q = \{o_i, \bullet_i \mid 1 \leq i \leq n\} \uplus \{q_f\}$  and we define the transition relation  $\Delta$  to contain the following transitions. The initial and final transitions are used to navigate from the root of the tree to the node matched by  $\sigma_1$ , and to read the final node matched by  $\sigma_n$  (and the selector as a whole) respectively. That is,  $o_1 \xrightarrow[\ast]{\downarrow, \rightarrow +} o_1$  and  $o_n \xrightarrow[\ast]{\circ} q_f$ . We have further transitions for  $1 \leq i < n$  that are shown in Figure 9. The transitions connect  $o_i$  to  $o_{i+1}$  depending on  $o_i$ . Figure 9a shows the child operator. The automaton moves to the first child of the current node and moves to the right zero or more steps. Figure 9b shows the descendant operator. The automaton traverses the tree downward and rightward any number of steps. The neighbour operator is handled in Figure 9c by simply moving to the next sibling. Finally, the sibling operator is shown in Figure 9d.

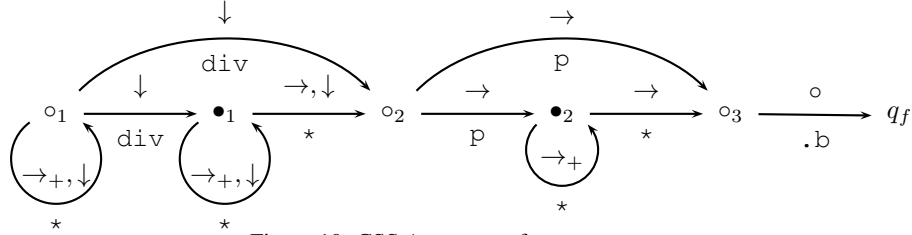
We prove the correctness of this construction in Lemma C.1 (soundness) and Lemma C.2 (completeness) in Appendix C.1.

### 6.2.3 Complex Example

Figure 10 gives an example of a CSS automaton representing the more complex selector

$$\text{div} \gg p \sim . b$$

which selects a node that has class  $b$ , is a right sibling of a node with element  $p$  and moreover is a descendent of a  $\text{div}$  node. This automaton again begins at  $o_1$  and navigates until it finds the node with the  $\text{div}$  element name. The automaton can read this node in two ways. The topmost transition covers

Figure 10: CSS Automaton for  $\text{div} \gg p \sim .b$ 

the case where the  $p$  node is directly below the  $\text{div}$  node. The lower transition allows the automaton to match the  $p$  node and then use a loop to navigate to the descendent node that will match  $p$ . Similarly, from  $o_2$  the automaton can read a  $p$  node and choose between immediately matching the node with class  $b$  or navigating across several siblings (using the loop at state  $\bullet_2$ ) before matching  $.b$  and accepting.

### 6.3 Closure Under Intersection

The problems of non-emptiness and intersection of CSS automata can be defined in precisely the same way we defined them for CSS selectors. One key property of CSS automata, which is not enjoyed by CSS selectors, is the closure of their languages under intersection. This allows us to treat the problem of intersection of CSS automata (i.e. the non-emptiness of the intersection of two CSS automata languages) as the non-emptiness problem (i.e. whether a given CSS automaton has an empty language).

**Proposition 6.2.** *Given two CSS automata  $\mathcal{A}_1$  and  $\mathcal{A}_2$ , we may construct in polynomial-time an automaton  $\mathcal{A}_1 \cap \mathcal{A}_2$  such that  $\mathcal{L}(\mathcal{A}_1) \cap \mathcal{L}(\mathcal{A}_2) = \mathcal{L}(\mathcal{A}_1 \cap \mathcal{A}_2)$ .*

The construction of the CSS automaton  $\mathcal{A}_1 \cap \mathcal{A}_2$  is by a variant of the standard product construction [71] for finite-state automata over finite words, which run the two given automata in parallel synchronised by the input word. Our construction runs the two CSS automata  $\mathcal{A}_1$  and  $\mathcal{A}_2$  in parallel synchronised by the path that they traverse. We first proceed with the formal definition and give an example afterwards.

#### 6.3.1 Formal Definition of Intersection

We first define the intersection of two node selectors. Recall node selectors are of the form  $\tau\Theta$  where  $\tau \in \{*, (s|*), e, (s|e) \mid s \in \text{NS} \wedge e \in \text{ELE}\}$ . The intersection of two node selectors  $\tau_1\Theta_1$  and  $\tau_2\Theta_2$  should enforce all properties defined in  $\Theta_1$  and  $\Theta_2$ . In addition, both selectors should be able to agree on the namespace and element name of the node, hence this part of the selector needs to be combined

more carefully. Thus, letting  $\Theta = \Theta_1 \cup \Theta_2$ , we define

$$\tau_1 \Theta_1 \cap \tau_2 \Theta_2 = \begin{cases} \tau_2 \Theta & \tau_1 = * \\ \tau_1 \Theta & \tau_2 = * \\ \tau_2 \Theta & \tau_1 = (s | *) \wedge \begin{pmatrix} \tau_2 = (s | *) \vee \\ \tau_2 = (s | e) \end{pmatrix} \\ (s | e) \Theta & \tau_1 = (s | *) \wedge \tau_2 = e \\ \tau_1 \Theta & \tau_1 = (s | e) \wedge \begin{pmatrix} \tau_2 = (s | e) \vee \\ \tau_2 = e \vee \\ \tau_2 = (s | *) \end{pmatrix} \\ \tau_2 \Theta & \tau_1 = e \wedge (\tau_2 = (s | e) \vee \tau_2 = e) \\ (s | e) \Theta & \tau_1 = e \wedge \tau_2 = (s | *) \\ \text{; not } (*) & \text{otherwise.} \end{cases}$$

We now define the automaton  $\mathcal{A}_1 \cap \mathcal{A}_2$ . The intersection automaton synchronises transitions that move in the same direction (by  $\downarrow$ ,  $\rightarrow$ ,  $\rightarrow_+$ ) or both agree to match the current node at the same time (with  $\circ$ ). In addition, we observe that a  $\rightarrow_+$  can be used by one automaton while the other uses  $\rightarrow$ . Given

$$\mathcal{A}_1 = (Q_1, \text{ELE}, \Delta_1, q_1^{\text{in}}, q_f^1) \quad \text{and} \quad \mathcal{A}_2 = (Q_2, \text{ELE}, \Delta_2, q_2^{\text{in}}, q_f^2)$$

we define

$$\mathcal{A}_1 \cap \mathcal{A}_2 = (Q_1 \times Q_2, \text{ELE}, \Delta, (q_1^{\text{in}}, q_2^{\text{in}}), (q_f^1, q_f^2))$$

where (letting  $d$  range over  $\{\rightarrow, \rightarrow_+, \downarrow, \circ\}$ ) we set  $\Delta =$

$$\begin{aligned} & \left\{ (q_1, q_2) \xrightarrow[\sigma_1 \cap \sigma_2]{d} (q'_1, q'_2) \mid q_1 \xrightarrow[\sigma_1]{d} q'_1 \wedge q_2 \xrightarrow[\sigma_2]{d} q'_2 \right\} \cup \\ & \left\{ (q_1, q_2) \xrightarrow[\sigma_1]{\rightarrow} (q'_1, q_2) \mid q_1 \xrightarrow[\sigma_1]{\rightarrow} q'_1 \wedge q_2 \xrightarrow[\ast]{\rightarrow_+} q_2 \right\} \cup \\ & \left\{ (q_1, q_2) \xrightarrow[\sigma_2]{\rightarrow} (q_1, q'_2) \mid q_1 \xrightarrow[\ast]{\rightarrow_+} q_1 \wedge q_2 \xrightarrow[\sigma_2]{\rightarrow} q'_2 \right\}. \end{aligned}$$

We verify that this transition relation satisfies the appropriate conditions:

1. (only self-loops) for a contradiction, a loop in  $\Delta$  that is not a self-loop can be projected to a loop of  $\mathcal{A}_1$  or  $\mathcal{A}_2$  that is also not a self-loop, e.g. if there exists  $(q_1, q_2) \xrightarrow[\sigma]{d} (q'_1, q'_2) \xrightarrow[\sigma']{d'} (q_1, q_2)$  that is not a self-loop, then either  $q_1 \neq q'_1$  or  $q_2 \neq q'_2$ , and thus we have a loop from  $q_1$  to  $q'_1$  to  $q_1$  in  $\mathcal{A}_1$  or similarly for  $\mathcal{A}_2$ ,
2. ( $\rightarrow_+$  loops and doesn't check nodes)  $\rightarrow_+$  transitions are built from  $\rightarrow_+$  transitions in  $\mathcal{A}_1$  and  $\mathcal{A}_2$ , thus a violation in the intersection implies a violation in one of the underlying automata,
3. ( $\rightarrow$  doesn't label loops)  $\rightarrow$  transitions are built from at least one  $\rightarrow$  transition in  $\mathcal{A}_1$  or  $\mathcal{A}_2$ , thus a violation in the intersection implies a violation in one of the underlying automata,
4. ( $\circ$  checks last node only) similarly, a violation of this constraint in the intersection implies a violation in one of the underlying automata,
5. ( $q_f$  is a sink) again, a violation of this constraint in the intersection implies a violation in the underlying automata.

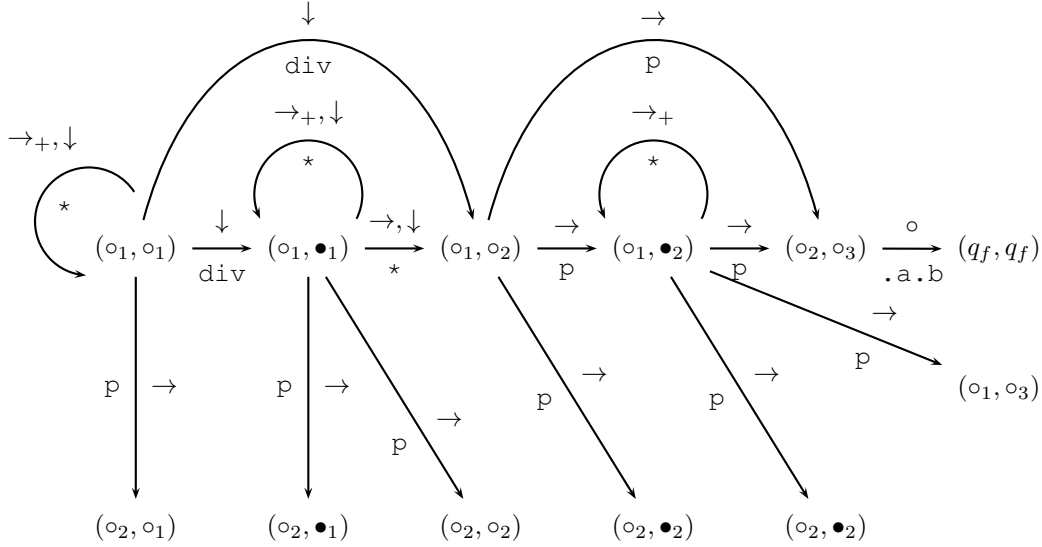


Figure 11: The intersection of the automaton in Figure 8 (in the first component) and the automaton in Figure 10 (in the second component).

### 6.3.2 Example of Intersection

Recall the automaton in Figure 8 (equivalent to  $p + \cdot a$ ) and the automaton in Figure 10 (equivalent to  $div \gg p \sim \cdot b$ ). The intersection of the two automata is given in Figure 11. Each state is a tuple  $(q_1, q_2)$  where the first component  $q_1$  represents the state of the automaton equivalent to  $p + \cdot a$  and the second component  $q_2$  the automaton equivalent to  $div \gg p \sim \cdot b$ .

In this example, accepting runs of the automaton will use only the top row of states. The lower states are reached when the two automata move out of sync and can no longer reach agreement on the final node matched. Usually this is by the first automaton matching a node labelled  $p$ , after which it must immediately accept the neighbouring node. This leaves the second automaton unable find a match. Hence, the first automaton needs to stay in state  $o_1$  until the second has reached a near-final state. Note, the two automata need not match the same node with element name  $p$ .

## 6.4 Reducing Non-emptiness of CSS Automata to SMT-solving

We will now provide a polynomial-time reduction from the non-emptiness of a CSS automaton to satisfiability of quantifier-free theory over integer linear arithmetic. That is, given a CSS automaton  $\mathcal{A}$ , our algorithm constructs a quantifier-free formula  $\theta_{\mathcal{A}}$  over integer linear arithmetic such that  $\mathcal{A}$  recognises a non-empty language iff  $\theta_{\mathcal{A}}$  is satisfiable. The encoding is quite involved and requires three small model properties discussed earlier. Once we have these properties we can construct the required formula of the quantifier-free theory over linear arithmetic. We begin by discussing each of these properties in turn, and then provide the reduction. The reduction is presented in a number of stages. We show how to handle attribute selectors separately before handling the encoding of CSS automata. The encoding of CSS automata is further broken down: we first describe the variables used in the encoding, then we

describe how to handle node selectors, finally we put it all together to encode runs of a CSS automaton.

For the remainder of the section, we fix a CSS automaton  $\mathcal{A} = (Q, \Delta, q^{\text{in}}, q_f)$  and show how to construct the formula  $\theta_{\mathcal{A}}$  in polynomial-time.

### 6.4.1 Bounded Run Length

The first property required is that the length of runs can be bounded. That is, if the language of  $\mathcal{A}$  is not empty, there is an accepting run over some tree whose length is smaller than the derived bound. We will construct a formula that will capture all runs of length up to this bound. Thanks to the bound we know that if an accepting run exists, the formula will encode at least one, and hence be satisfiable.

**Proposition 6.3** (Bounded Runs). *Given a CSS Automaton  $\mathcal{A} = (Q, \Delta, q^{\text{in}}, q_f)$ , if  $\mathcal{L}(\mathcal{A}) \neq \emptyset$ , there exists  $(T, \eta) \in \mathcal{L}(\mathcal{A})$  with an accepting run of length  $|\Delta|$ .*

This proposition is straightforward to obtain. We exploit that any loop in the automaton is a self loop that only needs to be taken at most once. For loops labelled  $\downarrow$ , a CSS formula cannot track the depth in the tree, so repeated uses of the loop will only introduce redundant nodes. For loops labelled  $\rightarrow_+$ , selectors such as `:nth-child( $\alpha n + \beta$ )` may enforce the existence of a number of intermediate nodes. But, since  $\rightarrow_+$  can cross several nodes, such loops also only needs to be taken once. Hence, each transition only needs to appear once in an accepting run. That is, if there is an accepting run of a CSS automaton with  $n$  transitions, there is also an accepting run of length at most  $n$ .

### 6.4.2 Bounding Namespaces and Elements

It will also be necessary for us to argue that the number of namespaces and elements can be bounded linearly in the size of the automaton. This is because our formula will need keep track of the number of nodes of each type appearing in the tree. This is required for encoding, e.g., the pseudo-classes of the form `:nth-of-type( $\alpha n + \beta$ )`. By bounding the number of types, our formula can use a bounded number of variables to store this information.

We state the property below. The proof is straightforward and appears in Appendix C.3.1. Intuitively, since only a finite number of nodes can be directly inspected by a CSS automaton, all others can be relabelled to a dummy type unless their type matches one of the inspected nodes.

**Proposition 6.4** (Bounded Types). *Given a CSS Automaton  $\mathcal{A} = (Q, \Delta, q^{\text{in}}, q_f)$  if there exists  $(T, \eta) \in \mathcal{L}(\mathcal{A})$  with  $T = (D, \lambda)$ , then there exists some  $(T', \eta) \in \mathcal{L}(\mathcal{A})$  where  $T' = (D, \lambda')$ . Moreover, let  $\downarrow(\text{ELE})$  be the set of element names and  $\downarrow(\text{NS})$  be the set of namespaces appearing in the image of  $\lambda'$ . Both the size of  $\downarrow(\text{ELE})$  and the size of  $\downarrow(\text{NS})$  are bounded linearly in the size of  $\mathcal{A}$ .*

### 6.4.3 Bounding Attribute Values

We will need to encode the satisfiability of conjunctions of attribute selectors. This is another potential source of unboundedness because the values are strings of arbitrary length. We show that, in fact, if the language of the automaton is not empty, there there is a solution whose attribute values are strings of length less than a bound polynomial in the size of the automaton.

The proof of the following lemma is highly non-trivial and uses techniques inspired by results in Linear Temporal Logic and automata theory. To preserve the flow of the article, we present the proof in Appendix C.3.2.

**Proposition 6.5** (Bounded Attributes). *Given a CSS Automaton  $\mathcal{A} = (Q, \Delta, q^{\text{in}}, q_f)$  if there exists  $(T, \eta) \in \mathcal{L}(\mathcal{A})$  with  $T = (D, \lambda)$ , then there exists some bound  $N$  polynomial in the size of  $\mathcal{A}$  and some  $(T', \eta) \in \mathcal{L}(\mathcal{A})$  where the length of all attribute values in  $T'$  is bound by  $N$ .*

Given such a bound on the length of values, we can use quantifier-free Presburger formulas to “guess” these witnessing strings by using a variable for each character position in the string. Then, the letter in each position is encoded by a number. This process is discussed in the next section.

#### 6.4.4 Encoding Attribute Selectors

Before discussing the full encoding, we first show how our formula can encode attribute selectors. Once we have this encoding, we can invoke it as a sub-routine of our main encoding whenever we have to handle attribute selectors. It is useful for readability reasons to present this in its own section.

The first key observation is that we can assume each positive attribute selector that does not specify a namespace applies to a unique, fresh, namespace. Thus, these selectors do not interact with any other positive attribute selectors and we can handle them easily. Note, these fresh namespaces do not appear in  $\downarrow(\text{NS})$ .

We present our encoding which works by identifying combinations of attribute selectors that must apply to the same attribute value. That is, we discover how many attribute values are needed, and collect together all selectors that apply to each selector. To that end, let  $op$  range over the set  $\{=, \sim, |=, \wedge, \$=, *=\}$  and let  $\tau\Theta$  be a node selector. For each  $s$  and  $a$ , let  $\Theta_a^s$  be the set of conditions in  $\Theta$  of the form  $\theta$  or  $:\text{not}(\theta)$  where  $\theta$  is of the form  $[s|a]$  or  $[s|a\ op\ v]$ . Recall we are encoding runs of a CSS automaton of length at most  $n$ . For a given position  $i$  in the run, we define  $\text{AttsPres}(\tau\Theta, i)$  to be the conjunction of the following constraints, where the encoding for  $\text{AttsPres}_{s:a}(\Theta, i)$  is presented below. Since a constraint of the form  $:\text{not}([a\ op\ v])$  applies to all attributes  $a$  regardless of their namespace, we define for convenience  $\text{Neg}(s, a) = \{:\text{not}([s|a\ op\ v]) \mid :\text{not}([a\ op\ v]) \in \Theta\}$ .

- For each  $s$  and  $a$  with  $\Theta_a^s$  non-empty and containing at least one selector of the form  $[s|a]$  or  $[s|a\ op\ v]$ , we enforce

$$\text{AttsPres}_{s:a}(\Theta_a^s \cup \text{Neg}(s, a), i)$$

if  $:\text{not}([a]) \notin \Theta$  and  $:\text{not}([s|a]) \notin \Theta$ , else we assert false.

- For each  $[a] \in \Theta$ , let  $s$  be fresh namespace. We assert

$$\text{AttsPres}_{s:a}(\{[s|a]\} \cup \text{Neg}(s, a), i)$$

and for each  $[a\ op\ v] \in \Theta$  we assert

$$\text{AttsPres}_{s:a}(\{[s|a\ op\ v]\} \cup \text{Neg}(s, a), i)$$

whenever, in both cases,  $:\text{not}([a]) \notin \Theta$ . If  $:\text{not}([a]) \in \Theta$  in both cases we assert false.

It remains to encode  $\text{AttsPres}_{s:a}(C, i)$  for some set of attribute selectors  $C$  all applying to  $s$  and  $a$ .

We can obtain a polynomially-sized global bound  $(N - 1)$  on the length of any satisfying value of an attribute  $s:a$  at some position  $i$  of the run from Proposition 6.5 (Bounded Attributes)<sup>6</sup>. Finally, we increment the bound by one to allow space for a trailing null character.

Once we have a bound on the length of a satisfying value, we can introduce variables  $x_{i,1}^{s:a}, \dots, x_{i,N}^{s:a}$  for each character position of the satisfying value, and encode the constraints almost directly. That is, letting  $\theta$  range over positive attribute selectors, we define<sup>7</sup>

$$\begin{aligned} \text{AttsPres}_{s:a}(C, i) = & \bigwedge_{\theta \in C} \text{AttsPres}(\theta, \vec{x}) \wedge \\ & \bigwedge_{:\text{not}(\theta) \in C} \neg \text{AttsPres}(\theta, \vec{x}) \wedge \text{Nulls}(\vec{x}). \end{aligned}$$

<sup>6</sup>Of course, we could obtain individual bounds for each  $s$  and  $a$  if we wanted to streamline the encoding.

<sup>7</sup>Note, we allow negation in this formula. This is for convenience only as the formulas we negate can easily be transformed into existential Presburger.

where  $\vec{x} = x_{i,1}^{s:a}, \dots, x_{i,N}^{s:a}$  will be existentially quantified later in the encoding and whose values will range<sup>8</sup> over  $\Gamma \uplus \{0\}$  where 0 is a null character used to pad the suffix of each word. We define  $\text{AttsPres}(\theta, \vec{x})$  for several  $\theta$ , the rest can be defined in the same way (see Appendix C.3.3). Letting  $v = a_1 \dots a_m$ ,

$$\begin{aligned} \text{AttsPres}([s | a], \vec{x}) &= \top \\ \text{AttsPres}([s | a = v], \vec{x}) &= \bigwedge_{1 \leq i \leq m} x_{i,j}^{s:a} = a_j \wedge x_{i,m+1}^{s:a} = 0 \\ \text{AttsPres}([s | a \wedge v], \vec{x}) &= \bigwedge_{1 \leq j \leq m} x_{i,j}^{s:a} = a_j \\ \text{AttsPres}([s | a * = v], \vec{x}) &= \bigvee_{0 \leq j \leq N-m-1} \bigwedge_{1 \leq j' \leq m} x_{i,j+j'}^{s:a} = a_j \end{aligned}$$

Finally, we enforce correct use of the null character

$$\text{Nulls}(\vec{x}) = \bigvee_{1 \leq j \leq N} \bigwedge_{j \leq j' \leq N} x_{i,j'}^{s:a} = 0.$$

#### 6.4.5 Encoding Non-Emptiness

We are now ready to give the main encoding of the emptiness of a CSS automaton using the quantifier-free theory over integer linear arithmetic. This encoding makes use of a number of variables, which we explain intuitively below. After describing the variables, we give the encoding in two parts: first we explain how a single node selector can be translated into existential Presburger arithmetic. Once we have this translation, we give the final step of encoding a complete run of an automaton.

**Variables Used in the Encoding** Our encoding makes use of the following variables for  $0 \leq i \leq n$ , representing the node at the  $i$ th step of the run. We use the overline notation to indicate variables.

- $\bar{q}_i$ , taking any value in  $Q$ , indicating the state of the automaton when reading the  $i$ th node in the run,
- $\bar{s}_i$ , taking any value in  $\downarrow(\text{NS})$  indicating the element tag (with namespace) of the  $i$ th node read in the run,
- $\bar{e}_i$ , taking any value in  $\downarrow(\text{ELE})$  indicating the element tag (with namespace) of the  $i$ th node read in the run,
- $\bar{p}_i$ , for each pseudo-class  $p \in P \setminus \{\text{root}\}$  indicating that the  $i$ th node has the pseudo-class  $p$ ,
- $\bar{n}_i$ , taking a natural number indicating that the  $i$ th node is the  $\bar{n}_i$ th child of its parent, and
- $\bar{n}_i^{s:e}$ , for all  $s \in \downarrow(\text{NS})$  and  $e \in \downarrow(\text{ELE})$ , taking a natural number variable indicating that there are  $\bar{n}_i^{s:e}$  nodes of type  $s:e$  strictly preceding the current node in the sibling order, and
- $\bar{N}_i$ , taking a natural number indicating that the current node is the  $\bar{N}_i$ th to last child of its parent, and
- $\bar{N}_i^{s:e}$ , for all  $s \in \downarrow(\text{NS})$  and  $e \in \downarrow(\text{ELE})$ , taking a natural number variable indicating that there are  $\bar{N}_i^{s:e}$  nodes of type  $s:e$  strictly following the current node in the sibling order, and

<sup>8</sup>Strictly speaking, Presburger variables range over natural numbers. It is straightforward to range over a finite number of values. That is, we can assume, w.l.o.g. that  $\Gamma \uplus 0 \subseteq \mathbb{N}$  and the quantification is suitably restricted.

- $x_{i,j}^{s:a}$  as used in the previous section for encoding the character at the  $j$ th character position of the attribute value for  $s:a$  at position  $i$  in the run<sup>9</sup>.

Note, we do not need a variable for `:root` since it necessarily holds uniquely at the 0th position of the run.

**Encoding Node Selectors** We define the encoding of node selectors below using the variables defined in the previous section. Note, this translation is not correct in isolation: global constraints such as “no ID appears twice in the tree” will be enforced later. The encoding works by translating each part of the selector directly. For example, the constraint  $e$  simply checks that  $\bar{e}_i = e$ . Even in the more complex cases of selectors such as `:nth-child( $\alpha n + \beta$ )` we are able to use a rather direct translation of the semantics:  $\exists \bar{n}. \bar{x} = \alpha \bar{n} + \beta$ . For the case of `:nth-of-type( $\alpha n + \beta$ )` we have to consider all possible namespaces  $s$  and element names  $e$  that the node could take, and use the  $\bar{n}_i^{s:e}$  variables to do the required counting.

In our presentation we allow ourselves to negate existentially quantified formulas of the form  $\exists \bar{n}. \bar{x} = \alpha \bar{n} + \beta$  where  $\bar{x}$  is a variable, and  $\alpha$  and  $\beta$  are constants. Although this is not strictly allowed in existential Presburger arithmetic, it is not difficult to encode correctly. For completeness, we provide the encoding of such negated formulas in Appendix C.3.4.

In the following, let  $\text{NoAtts}(\Theta)$  be  $\Theta$  less all selectors of the form `[ $s|a$ ]`, `[ $s|a$   $op$   $v$ ]`, `[ $a$ ]`, or `[ $a$   $op$   $v$ ]`, or `:not([ $s|a$ ])`, `:not([ $s|a$   $op$   $v$ ])`, `:not([ $a$ ])`, or `:not([ $a$   $op$   $v$ ])`.

**Definition 6.6** ( $\text{Pres}(\sigma, i)$ ). *Given a node selector  $\tau\Theta$ , we define*

$$\text{Pres}(\tau\Theta, i) = \left( \begin{array}{c} \text{Pres}(\tau, i) \wedge \\ \left( \bigwedge_{\theta \in \text{NoAtts}(\Theta)} \text{Pres}(\theta, i) \right) \wedge \\ \text{AttsPres}(\tau\Theta, i) \end{array} \right)$$

where we define  $\text{Pres}(\theta, i)$  as follows:

$$\begin{aligned} \text{Pres}(*, i) &= \top \\ \text{Pres}((s|*), i) &= (\bar{s}_i = s) \\ \text{Pres}(e, i) &= (\bar{e}_i = e) \\ \text{Pres}((s|e), i) &= (\bar{s}_i = s \wedge \bar{e}_i = e) \\ \text{Pres}(:\text{not}(\sigma\_), i) &= \neg \text{Pres}(\sigma\_, i) \\ \text{Pres}(:\text{root}, i) &= \begin{cases} \top & i = 0 \\ \perp & \text{otherwise} \end{cases} \\ \forall p \in P \setminus \{:\text{root}\}. \text{Pres}(p, i) &= \bar{p}_i \end{aligned}$$

and, finally, for the remaining selectors, we have

$$\begin{aligned} \text{Pres}(:\text{nth-child}(\alpha n + \beta), 0) &= \perp \\ \text{Pres}(:\text{nth-last-child}(\alpha n + \beta), 0) &= \perp \\ \text{Pres}(:\text{nth-of-type}(\alpha n + \beta), 0) &= \perp \\ \text{Pres}(:\text{nth-last-of-type}(\alpha n + \beta), 0) &= \perp \\ \text{Pres}(:\text{only-child}, 0) &= \perp \\ \text{Pres}(:\text{only-of-type}, 0) &= \perp \end{aligned}$$

<sup>9</sup>Recall  $s$  is not necessarily in  $\downarrow(\text{NS})$  as it may be some fresh value.



and when  $i > 0$

$$\begin{aligned}
\text{Pres}(\text{:nth-child}(\alpha n + \beta), i) &= \exists \bar{n}. \bar{n}_i = \alpha \bar{n} + \beta \\
\text{Pres}(\text{:nth-last-child}(\alpha n + \beta), i) &= \exists \bar{n}. \bar{N}_i = \alpha \bar{n} + \beta \\
\text{Pres}(\text{:nth-of-type}(\alpha n + \beta), i) &= \bigvee_{\substack{s \in \downarrow(\text{NS}) \\ e \in \downarrow(\text{ELE})}} \left( \begin{array}{l} \bar{s}_i = s \wedge \bar{e}_i = e \wedge \\ \exists \bar{n}. \bar{n}_i^{s:e} + 1 = \alpha \bar{n} + \beta \end{array} \right) \\
\text{Pres}(\text{:nth-last-of-type}(\alpha n + \beta), i) &= \bigvee_{\substack{s \in \downarrow(\text{NS}) \\ e \in \downarrow(\text{ELE})}} \left( \begin{array}{l} \bar{s}_i = s \wedge \bar{e}_i = e \wedge \\ \exists \bar{n}. \bar{N}_i^{s:e} + 1 = \alpha \bar{n} + \beta \end{array} \right) \\
\text{Pres}(\text{:only-child}, i) &= \bar{n}_i = 1 \wedge \bar{N}_i = 1 \\
\text{Pres}(\text{:only-of-type}, i) &= \bigvee_{\substack{s \in \downarrow(\text{NS}) \\ e \in \downarrow(\text{ELE})}} \left( \begin{array}{l} \bar{s}_i = s \wedge \bar{e}_i = e \wedge \\ \bar{n}_i^{s:e} = 0 \wedge \bar{N}_i^{s:e} = 0 \end{array} \right)
\end{aligned}$$

We are now ready to move on to complete the encoding.

**Encoding Runs of CSS Automata** Finally, now that we are able to encode attribute and node selectors, we can make use of these to encode accepting runs of a CSS automaton. Since we know that, if there is an accepting run, then there is a run of length at most  $n$  where  $n$  is the number of transitions in  $\Delta$ , we encode the possibility of an accepting run using the variables discussed above for all  $0 \leq i \leq n$ . The shape of the translation is given below and elaborated on afterwards.

**Definition 6.7.**  $\theta_{\mathcal{A}}$  Given a CSS automaton  $\mathcal{A}$  we define

$$\theta_{\mathcal{A}} = \left( \left( \begin{array}{c} \bar{q}_0 = q^{in} \\ \wedge \\ \bar{q}_n = q_f \end{array} \right) \wedge \bigwedge_{0 \leq i < n} \left( \begin{array}{c} \text{Tran}(i) \\ \vee \\ \bar{q}_i = q_f \end{array} \right) \wedge \text{Consistent} \right)$$

where  $\text{Tran}(i)$  and  $\text{Consistent}$  are defined below.

Intuitively, the first two conjuncts asserts that a final state is reached from an initial state. Next, we use  $\text{Tran}(i)$  to encode a single step of the transition relation, or allows the run to finish early. Finally  $\text{Consistent}$  asserts consistency constraints.

We define as a disjunction over all possible (single-step) transitions  $\text{Tran}(i) = \bigvee_{t \in \Delta} \text{Tran}(i, t)$  where  $\text{Tran}(i, t)$  is defined below by cases. There are four cases depending on whether the transition is labelled  $\downarrow$ ,  $\rightarrow$ ,  $\rightarrow_+$ , or  $\circ$ . In most cases, we simply assert that the state changes as required by the transition, and that the variables  $\bar{n}_i$  and  $\bar{n}_i^{s:e}$  are updated consistently with the number of nodes read by the transition. Although the encodings look complex, they are essentially simple bookkeeping.

To ease presentation, we write  $\bar{s}_i : \bar{e}_i = s:e$  as shorthand for  $(\bar{s}_i = s \wedge \bar{e}_i = e)$  and  $\bar{s}_i : \bar{e}_i \neq s:e$  as shorthand for  $(\bar{s}_i \neq s \vee \bar{e}_i \neq e)$ .

1. When  $t = q \xrightarrow[\sigma]{\downarrow} q'$  we define  $\text{Tran}(i, t)$  to be

$$\begin{aligned}
&(\bar{q}_i = q) \wedge (\bar{q}_{i+1} = q') \wedge \neg \overline{\text{empty}_i} \wedge \text{Pres}(\sigma, i) \wedge \\
&(\bar{n}_{i+1} = 1) \wedge \bigwedge_{\substack{s \in \downarrow(\text{NS}) \\ e \in \downarrow(\text{ELE})}} (\bar{n}_{i+1}^{s:e} = 0) .
\end{aligned}$$

2. When  $t = q \xrightarrow[\sigma]{\rightarrow} q'$  we define  $\text{Tran}(i, t)$  to be false when  $i = 0$  (since the root has no siblings) and otherwise

$$\bigwedge_{\substack{s \in \downarrow(\text{NS}) \\ e \in \downarrow(\text{ELE})}} \left( \begin{array}{l} (\bar{q}_i = q) \wedge (\bar{q}_{i+1} = q') \wedge \text{Pres}(\sigma, i) \wedge \\ (\bar{n}_{i+1} = \bar{n}_i + 1) \wedge (\bar{N}_{i+1} = \bar{N}_i - 1) \wedge \\ \left( \begin{array}{l} (\bar{s}_i : \bar{e}_i = s:e) \Rightarrow (\bar{n}_{i+1}^{s:e} = \bar{n}_i^{s:e} + 1) \\ (\bar{s}_i : \bar{e}_i \neq s:e) \Rightarrow (\bar{n}_{i+1}^{s:e} = \bar{n}_i^{s:e}) \end{array} \right) \wedge \\ \left( \begin{array}{l} (\bar{s}_{i+1} : \bar{e}_{i+1} = s:e) \Rightarrow (\bar{N}_{i+1}^{s:e} = \bar{N}_i^{s:e} - 1) \\ (\bar{s}_{i+1} : \bar{e}_{i+1} \neq s:e) \Rightarrow (\bar{N}_{i+1}^{s:e} = \bar{N}_i^{s:e}) \end{array} \right) \end{array} \right) \wedge \end{array} \right).$$

3. When  $t = q \xrightarrow[*]{\rightarrow+} q$  we define  $\text{Tran}(i, t)$  to be false when  $i = 0$  and otherwise

$$\bigwedge_{\substack{s \in \downarrow(\text{NS}) \\ e \in \downarrow(\text{ELE})}} \exists \bar{\delta}_{s:e} \cdot \left( \begin{array}{l} (\bar{q}_i = q) \wedge (\bar{q}_{i+1} = q) \wedge \\ \exists \bar{\delta}. \left( (\bar{n}_{i+1} = \bar{n}_i + \bar{\delta}) \wedge (\bar{N}_{i+1} = \bar{N}_i - \bar{\delta}) \right) \wedge \\ \left( \begin{array}{l} (\bar{s}_i : \bar{e}_i = s:e) \Rightarrow \\ (\bar{n}_{i+1}^{s:e} = \bar{n}_i^{s:e} + \bar{\delta}_{s:e} + 1) \\ (\bar{s}_i : \bar{e}_i \neq s:e) \Rightarrow \\ (\bar{n}_{i+1}^{s:e} = \bar{n}_i^{s:e} + \bar{\delta}_{s:e}) \end{array} \right) \wedge \\ \left( \begin{array}{l} (\bar{s}_{i+1} : \bar{e}_{i+1} = s:e) \Rightarrow \\ (\bar{N}_{i+1}^{s:e} = \bar{N}_i^{s:e} - \bar{\delta}_{s:e} - 1) \\ (\bar{s}_{i+1} : \bar{e}_{i+1} \neq s:e) \Rightarrow \\ (\bar{N}_{i+1}^{s:e} = \bar{N}_i^{s:e} - \bar{\delta}_{s:e}) \end{array} \right) \end{array} \right) \wedge \end{array} \right).$$

4. When  $t = q \xrightarrow[\sigma]{\circ} q'$  we define  $\text{Tran}(i, t)$  to be

$$(\bar{q}_i = q) \wedge (\bar{q}_{i+1} = q') \wedge \text{Pres}(\sigma, i).$$

To ensure that the run is over a consistent tree, we assert the consistency constraint

$$\text{Consistent} = \text{Consistent}_n \wedge \text{Consistent}_i \wedge \text{Consistent}_p$$

where each conjunct is defined below.

- The clause  $\text{Consistent}_n$  asserts that the values of  $\bar{n}_i$ ,  $\bar{N}_i$ ,  $\bar{n}_i^{s:e}$ , and  $\bar{N}_i^{s:e}$  are consistent. That is

$$\bigwedge_{1 \leq i \leq n} \left( \bar{n}_i = 1 + \sum_{s:e \in \text{ELE}} \bar{n}_i^{s:e} \right) \wedge \left( \bar{N}_i = 1 + \sum_{s:e \in \text{ELE}} \bar{N}_i^{s:e} \right).$$

- The clause  $\text{Consistent}_i$  asserts that ID values are unique. It is the conjunction of the following clauses. For each  $s$  for which we have created variables of the form  $x_{i,j}^{s:\text{id}}$  we assert

$$\bigwedge_{1 \leq i \neq i' \leq n} \bigvee_{1 \leq j \leq N} x_{i,j}^{s:\text{id}} \neq x_{i',j}^{s:\text{id}}.$$

- Finally,  $\text{Consistent}_p$  asserts the remaining consistency constraints on the pseudo-classes. We define  $\text{Consistent}_p =$

$$\bigwedge_{0 \leq i \leq n} \left( \begin{array}{c} \neg(\text{:link}_i \wedge \text{:visited}_i) \wedge \\ \bigwedge_{0 \leq j \neq i \leq n} (\neg(\text{:target}_i \wedge \text{:target}_j)) \wedge \\ \neg(\text{:enabled}_i \wedge \text{:disabled}_i) \end{array} \right).$$

These conditions assert the mutual exclusivity of `:link` and `:visited`, that at most one node in the document can be the target node, that nodes are not both enabled and disabled.

#### 6.4.6 Correctness of the Encoding

We have now completed the definition of the reduction from the emptiness problem of CSS automata to the satisfiability of existential Presburger arithmetic. What remains is to show that this reduction is faithful: that is, the CSS automaton has an empty language if and only if the formula is satisfiable. The proof is quite routine, and presented in Lemma C.8 and Lemma C.9 in Appendix C.3.5.

**Lemma 6.8** (Correctness of  $\theta_{\mathcal{A}}$ ). *For a CSS automaton  $\mathcal{A}$ , we have*

$$\mathcal{L}(\mathcal{A}) \neq \emptyset \Leftrightarrow \theta_{\mathcal{A}} \text{ is satisfiable.}$$

We are thus able to decide the emptiness problem, and therefore the emptiness of intersection problem, for CSS automata by reducing the problem to satisfiability of existential Presburger arithmetic and using a fast solver such as Z3 [18] to resolve the satisfiability.

## 7 Rule-Merging to Max-SAT

In this section, we provide a reduction from the rule-merging problem to partial weighted MaxSAT. The input will be a valid covering  $\mathcal{C} = \{\overline{B}_i\}_{i=1}^m$  of a CSS graph  $\mathcal{G}$ . We aim to find a rule  $\overline{B} = (X, \overline{Y})$  and a position  $j$  that minimises the weight of  $\mathcal{C}[\overline{B} \rightarrow j]_{\downarrow}$ .

There is a fairly straightforward encoding of the rule-merging problem into Max-SAT using for each node  $w \in SUP$  a boolean variable  $\overline{w}$  which is true iff the node is included in the new rule. Unfortunately, early experiments showed that such an encoding does not perform well in practice, causing prohibitively long run times even on small examples. The failure of the naive encoding may be due to the search space that includes a large number of possible pairs  $(X, \overline{Y})$  that turn out to be invalid rules (e.g. include edges not in the CSS-graph  $\mathcal{G}$ ). Hence, we will use a different Max-SAT encoding that, by means of syntax, further restricts the search space of valid rule-merging opportunities.

The crux of our new encoding is to explicitly say in the Max-SAT formula  $\varphi$  that the rule  $\overline{B}$  in a merging opportunity  $(\overline{B}, j)$  is a “valid” sub-biclique of one of the *maximal* bicliques  $B = (X, Y)$  — maximal with respect to subset-of relations of  $X$  and  $Y$  components of bicliques — in the CSS-graph  $\mathcal{G}$ . By insisting  $\overline{B}$  is contained within a maximal biclique of the CSS-graph, we automatically ensure that  $\overline{B}$  does not contain edges that are not in  $\mathcal{G}$ .

The formula  $\varphi$  will try to guess a maximal biclique  $B$  and which nodes to omit from  $B$ . Since the number of maximal bicliques in a bipartite graph is exponential (in the number of nodes) in the worst case, one concern with this idea is that the constraint  $\varphi$  might become prohibitively large. As we shall see, this turns out not to be the case in practice. Intuitively, based on our experience, the number of rules in a real-world CSS file is *at most* a few thousand. Second, the number of maximal bicliques in a CSS-graph that corresponds to a real-world CSS file also is typically of a similar size to the number of rules, and, furthermore, can be enumerated using the algorithm from [36] (which runs in time polynomial in

```

.a { color:blue; color:green }
.b { color:green; color:blue }
```

Figure 12: A CSS file with an unorderable sub-biclique.

the size of the input and the size of the output). To be more precise, the benchmarks in our experiments had between 31 and 2907 rules, and the mean number of rules was 730. The mean ratio of the number of maximal bicliques to the number of rules was 1.25, and the maximum was 2.05. As we shall see in Section 8, Z3 may solve the constraints via this encoding quite efficiently.

In the rest of the section, we will describe our encoding in detail. For convenience, our encoding also allows bounded integer variables. There are standard ways to encode these in binary as booleans (e.g. see [53]) by bit-blasting (using a logarithmic number of boolean variables).

## 7.1 Orderable Bicliques

Our description above of the crux of our encoding (by restricting to containment in a maximal biclique) is a simplification. This is because *not all* sub-bicliques of a maximal biclique correspond to a valid rule  $\bar{B}$  in a merging opportunity  $(\bar{B}, j)$  with respect to the covering  $\mathcal{C}$ . [A biclique  $(X', Y')$  is a *sub-biclique* of a biclique  $(X, Y)$  if  $X' \subseteq X$  and  $Y' \subseteq Y$ .] To ensure that our constraint  $\varphi$  chooses only valid rules, it needs to ensure that the sub-biclique that is chosen is “orderable”. More precisely, a biclique  $B = (X, Y)$  is *orderable at position  $j$*  if it can be turned into a rule  $\bar{B} = (X, \bar{Y})$  (i.e. turning the set  $Y$  of declarations into a sequence  $\bar{Y}$  by assigning an order) that can be inserted at position  $j$  in  $\mathcal{C}$  without violating the validity of the resulting covering with respect to the order  $\prec$  (from the CSS-graph  $\mathcal{G}$ ). If there are  $m$  rules in  $\mathcal{C}$ , there are  $m + 1$  positions (call these positions  $0, \dots, m$ ) where  $\bar{Y}$  may be inserted into  $\mathcal{C}$ . We show below that the position  $j$  is crucial to whether  $\bar{B}$  is orderable.

Unorderable bicliques rarely arise in practice (in our benchmarks, the mean percentage of maximal bicliques that were unorderable at some position was 0.39%), but they have to be accounted for if our analysis is to find the optimal rule-merging opportunity. A biclique  $B = (X, Y)$  is unorderable when they have the same property name (or two related property names, e.g., shorthands) occurring multiple times with different values in  $Y$ . One reason having a CSS rule with the same property name occurring multiple times with different values is to provide “fallback options” especially because old browsers may not support certain values in some property names, e.g., the rule

```
.c { color:#ccc; color:rgba(0, 0, 0, 0.5); }
```

says that if `rgba()` is not supported (e.g. in IE8 or older browsers), then use `#ccc`. Using this idea, we can construct the simple example of an unorderable biclique in the CSS file in Figure 12. The ordering constraints we can derive from this file include

$$(.a, \text{color:blue}) \prec (.a, \text{color:green})$$

and

$$(.b, \text{color:green}) \prec (.b, \text{color:blue}).$$

The biclique  $B$

$$(\{.a, .b\}, \{\text{color:blue, color:green}\})$$

is easily seen to be orderable at position 0 and 1. This is because the final rule in Figure 12 will ensure the ordering  $(.b, \text{color:green}) \prec (.b, \text{color:blue})$  is satisfied, and since  $B$  will appear before this final rule, only  $(.a, \text{color:blue}) \prec (.a, \text{color:green})$  needs to be maintained by  $B$  (in fact, at position 0, neither of the orderings need to be respected by  $B$ ). However, at position 2, which is at the

end of the file in Figure 12, both orderings will have to be respected by  $B$ . Unfortunately, one of these orderings will be violated regardless of how one may assign an ordering to `blue` and `green`. This contrived example was made only for illustration, however, our technique should still be able to handle even contrived examples.

We mention that both checking orderability and ordering a given biclique can be done efficiently.

**Proposition 7.1.** *Given a biclique  $B$ , a covering  $\mathcal{C}$  (with  $m$  rules) of a CSS-graph  $\mathcal{G}$ , and a number  $j \in \{0, \dots, m\}$ , checking whether  $B$  is orderable at position  $j$  in  $\mathcal{C}$  can be done in polynomial time. Moreover, if  $B$  is orderable an ordering can be calculated in polynomial time.*

The proof of the proposition is easy and is relegated into Appendix A.1.

### Maximal Orderable Bicliques

Our Max-SAT encoding  $\varphi$  needs to ensure that we only pick a pair  $(B, j)$  such that  $B$  is an orderable biclique at position  $j$  in the given covering  $\mathcal{C}$ , i.e.,  $B$  corresponds to a rule that can be inserted at position  $j$  in  $\mathcal{C}$ . Although the check of orderability can be *declaratively* expressed as a constraint in  $\varphi$ , we found that this results in Max-SAT formulas that are rather difficult to solve by existing Max-SAT solvers. For this reason, we propose to express the check of orderability in a different way. Intuitively, for each  $j \in \{0, \dots, m\}$ , we enumerate all orderable bicliques  $B = (X, Y)$  that are also maximal, i.e., it is *not* a (strict) sub-biclique of a different orderable biclique. Since “orderability is inherited by sub-bicliques” (as the following lemma, whose proof is immediate from the definition, states), the constraint  $\varphi$  needs to simply choose a sub-biclique of a maximal orderable biclique that appears in our enumeration.

**Lemma 7.2.** *Every sub-biclique  $B' = (X, Y)$  of an orderable biclique  $B = (X, Y)$  is orderable.*

The above enumeration of maximal orderable bicliques can be described as a pair  $(\{M_i\}_{i=1}^\mu, \mathcal{F})$  where

- $\{M_i\}_{i=1}^\mu$  is an enumeration of all bicliques that are orderable and maximal at some position  $j$ , and
- $\mathcal{F}$  *forbids* certain bicliques at each position. I.e. it is a function from  $[1, m]$  to the set of bicliques in  $\{M_i\}_{i=1}^\mu$  that are unorderable at position  $j$ .

Observe that the set of orderable bicliques at position  $j$  in  $\mathcal{C}$  is a subset of the set of orderable bicliques at position  $j + 1$  in  $\mathcal{C}$ . This may be formally expressed as:  $\mathcal{F}(j) \subseteq \mathcal{F}(j + 1)$  for all  $j \in [1, m)$ .

In the majority (54%) of examples that we have from real-world CSS, the function  $\mathcal{F}$  maps all values of  $[1, m]$  to  $\emptyset$ , i.e., all maximal bicliques are orderable at all positions. The mean percentage of maximal bicliques that were unorderable at some position was 0.39%, with a maximum of 5.84%.

In our description of the Max-SAT encoding below, we assume that the pair  $(\{M_i\}_{i=1}^\mu, \mathcal{F})$  has been computed for the input  $\mathcal{C}$ .

## 7.2 The Max-SAT Encoding

We present the full reduction of the rule-merging problem to Max-SAT. In particular, the constraints we produce are

$$(\Pi_H, \Pi_S)$$

where  $\Pi_H$  and  $\Pi_S$  are, respectively, hard and soft constraints. First, we describe the variables used in our encoding. Note, our encoding will rely on the assumption that covering  $\mathcal{C}$  has already been trimmed. Recall, the notion of trimming is defined in Section 4 and is the process of removing redundant nodes from a covering. A node is redundant in a rule if all of its incident edges also appear later in the covering.

### 7.2.1 Representing the rule-merging opportunity

We need to represent a merging opportunity  $(\bar{B}, j)$ . We use a bounded integer variable  $\bar{j}$  (with range  $0 \leq \bar{j} \leq m$ ) to encode  $j$ .

For  $\bar{B}$  we select a biclique in  $\{M_i\}_{i=1}^\mu$  and allow some nodes to be removed (i.e. to produce sub-bicliques of the  $M_i$ ). We use a bounded integer variable  $\bar{i}_M$  (with range  $[1, \mu]$ ) to select  $M_i$ . Next, we need to choose a sub-biclique of  $M_i$ , which can be achieved by choosing nodes  $M_i$  to be removed. To minimise the number of variables used, we number the nodes contained in each biclique in some (arbitrary) way, i.e., for each  $i \in [1, \mu]$  and biclique  $M_i = (X, Y)$ , we define a bijection  $\rho_i : X \cup Y \rightarrow [1, |X \cup Y|]$ . Let  $\chi$  be the maximum number of nodes in a biclique  $M_i$  in the enumeration  $\{M_i\}_{i=1}^\mu$ , i.e., the maximal integer  $k$  such that  $\rho_i(w) = k$  for some  $w \in S \cup P$  and  $1 \leq i \leq \mu$ .

We introduce boolean variables  $\bar{x}_1, \dots, \bar{x}_\chi$ . Once the maximal orderable biclique  $M_i$  is picked, for a node  $w$  with  $\rho_i(w) = k$ , the variable  $\bar{x}_k$  is used to indicate that  $w$  is to be excluded from selected  $M_i$  (i.e.  $\bar{x}_k$  is true iff  $w$  is to be excluded from  $M_i$ ). More precisely, for an edge  $e = (s, p)$  in  $M_i$ , we define a predicate  $\text{HasEdge}(e)$  to be

$$\text{HasEdge}((s, p)) = \bigvee_{1 \leq i \leq \mu} \bar{i}_M = i \wedge \neg \bar{x}_{\rho_i(s)} \wedge \neg \bar{x}_{\rho_i(p)}.$$

Note, when  $M_i = (X, Y)$  and  $w \notin X \cup Y$  we let  $\bar{x}_{\rho_i(w)}$  denote the formula “true”.

### 7.2.2 Hard Constraints

We define

$$\Pi_H = \{\varphi_{\text{vld}}, \varphi_{\text{ord}}\}$$

where  $\varphi_{\text{vld}}$  and  $\varphi_{\text{ord}}$  are described below.

We need to ensure the rule-merging opportunity is valid, i.e., it has not been forbidden at the chosen position and inserting it into the covering does not violate the edge order  $\prec$ . For the former, we define  $\mathcal{F}^{\text{1st}}$ , which is used to discover which of the  $M_i$  first become unorderable at position  $j$ . That is  $M_i \in \mathcal{F}^{\text{1st}}(j)$  if  $j$  is the smallest integer such that  $M_i \in \mathcal{F}(j)$ .

$$\varphi_{\text{vld}} = \bigwedge_{1 \leq j \leq m} \left( (\bar{j} \geq j) \Rightarrow \bigwedge_{M_i \in \mathcal{F}^{\text{1st}}(j)} (\bar{i}_M \neq i) \right).$$

We also need to ensure the edge ordering is respected by the rule-merging opportunity, for which we define  $\varphi_{\text{ord}}$ . If  $e_1 \prec e_2$ , and  $e_1 = (s_1, p_1)$  is in the rule  $\bar{B}$  in the guessed merging opportunity  $(\bar{B}, j)$ , then we need to assert  $e_1 \prec e_2$  is respected. It is respected either if  $e_2 = (s_2, p_2)$  appears after the position where  $j$  is to be inserted in  $\mathcal{C}$ , or  $e_2$  is also contained in the new rule  $\bar{B}$  (in which case the ordering can still be respected since  $\bar{B}$  is orderable). That is,

$$\varphi_{\text{ord}} = \left( \bigwedge_{(s_1, p_1) \prec (s_2, p_2)} \text{HasEdge}((s_1, p_1)) \Rightarrow (\bar{j} < \text{index}((s_2, p_2)) \vee \text{HasEdge}((s_2, p_2))) \right).$$

This is because only the last occurrence of an edge in a covering matters (recall the definition of index of an edge in Section 4). Also note that our use of bicliques ensures that we do not introduce pairs  $(s, p)$  that are not edges in  $E$ .

### 7.2.3 Soft Constraints

The soft constraints will calculate the weight of  $\mathcal{C}[\overline{B} \rightarrow j]_{\downarrow}$ . Since we want to minimise this weight, the optimal solution to the Max-SAT problem will give an optimal rule-merging opportunity. Our soft constraints are

$$\Pi_S = \Pi_{\overline{S}} \cup \Pi_S^{\text{Sels}} \cup \Pi_S^{\text{Props}}$$

where  $\Pi_{\overline{S}}$  counts the weight of the  $\overline{B}$ , and  $\Pi_S^{\text{Sels}}$  and  $\Pi_S^{\text{Props}}$  counts the weight of the remainder (not including  $\overline{B}$ ) of the stylesheet by counting the remaining selectors and properties respectively after trimming. These are defined below.

To count the weight of  $\overline{B}$ , we count the weight of the non-excluded nodes of  $M_i = (X_i, Y_i)$ . We have

$$\Pi_{\overline{S}} = \{(\varphi_w^i, \text{wt}(w)) \mid 1 \leq i \leq \mu \wedge w \in X_i \cup Y_i\}$$

where

$$\varphi_w^i = ((\overline{i}_M = i) \Rightarrow \overline{x}_{\rho_i(w)}) .$$

Note that  $\varphi_w^i$  is true iff, whenever  $M_i$  is picked, the node  $w$  is omitted from  $M_i$  in the guessed  $\overline{B}$ . Furthermore, the cost of *not* omitting  $w$  from  $M_i$  is  $\text{wt}(w)$ .

Next, we count the remaining weight of  $\mathcal{C}[\overline{B} \rightarrow j]$  (i.e. excluding the weight of  $\overline{B}$ ). Assume that  $\mathcal{C}$  has already been trimmed, i.e.,  $\mathcal{C}_{\downarrow} = \mathcal{C}$ . The intuition is that a node  $v$  can be removed from  $\overline{B}_i$  in  $\mathcal{C}$  if all edges  $e$  incident to  $v$  appear in a later rule in  $\mathcal{C}$ , i.e.,  $\text{index}(e) > i$ . In particular, let  $\overline{B}_i = (X_i, \overline{Y}_i)$  in  $\{\overline{B}_i\}_{i=1}^m$ . To count the weight of the untrimmed selectors we use the clauses

$$\Pi_S^{\text{Sels}} = \{(\psi_s^i, \text{wt}(s)) \mid 1 \leq i \leq m \wedge s \in X_i\}$$

where

$$\psi_s^i = \left( i \leq \overline{j} \wedge \bigwedge_{\substack{\text{index}((s,p))=i \\ p \in \overline{Y}}} \text{HasEdge}((s,p)) \right) .$$

The idea is that  $\psi_s^i$  will be satisfied whenever  $s$  can be removed from  $\overline{B}_i$ . We assume  $\mathcal{C}$  has already been trimmed, so  $s$  can only become removable because of the application of rule-merging. This explains the first conjunct which asserts that nodes can only be removed from rules appearing before  $\overline{j}$ . Next,  $s$  can be removed after rule-merging if none of its incident edges  $(s, p)$  are the final occurrence of  $(s, p)$  in  $\mathcal{C}[\overline{B} \rightarrow j]$ . The crucial edges in this check are those such that  $\text{index}((s, p)) = i$ , which means  $\overline{B}_i$  contains the final occurrence of  $(s, p)$  in  $\mathcal{C}$ . For these edges, if  $\text{HasEdge}((s, p))$  holds, then the new rule contains  $(s, p)$  and  $\overline{B}_i$  will no longer contain the final occurrence of  $(s, p)$  after rule-merging.

Similarly, to count the weight of the properties that cannot be removed we have

$$\Pi_S^{\text{Props}} = \{(\psi_p^i, \text{wt}(p)) \mid 1 \leq i \leq m \wedge p \in \overline{Y}_i\}$$

where

$$\psi_p^i = \left( i \leq \overline{j} \wedge \bigwedge_{\substack{\text{index}((s,p))=i \\ s \in X}} \text{HasEdge}((s,p)) \right) .$$

### 7.3 Generated Rule-Merging Opportunity

The merging opportunity  $(\bar{B}, j)$  is built from an optimal satisfying assignment to  $(\Pi_H, \Pi_S)$ . First,  $j$  is the value assigned to  $\bar{j}$ . Then  $\bar{B} = (X, \bar{Y})$  where, letting  $i_M$  be the value of  $\bar{i}_M$ , and letting  $M_{i_M} = (X', Y')$ ,

- $s \in X$  iff  $s \in X'$  and  $\bar{x}_{\rho_{i_M}(s)}$  is assigned the value false, and
- $\bar{Y} = \{p_i\}_{i=1}^m$ , where  $\{p_i\}_{i=1}^m$  is obtained by assigning an ordering to  $Y'$  such that  $(\bar{B}, j)$  is a valid merging opportunity.

That the ordering of  $Y'$  above exists is guaranteed by the fact that  $M_{i_M}$  is orderable at  $j$ , and Lemma 7.2. We can compute the ordering in polynomial time via Proposition 7.1. We argue the following proposition in Appendix A.3.

**Proposition 7.3.** *The merging opportunity  $(\bar{B}, j)$  generated from the maximal solution to  $(\Pi_H, \Pi_S)$  is the optimal merging opportunity of  $\mathcal{C}$ .*

## 8 Experimental Results

We implemented a tool SATCSS (in Python 2.7) for CSS minification via rule-merging which can be found in our supplementary material. The tool is also available as source code and in a working disk image via the following URLs.

<https://github.com/matthewhague/sat-css-tool>  
<http://www.cs.rhul.ac.uk/home/hague/sat-css-tool.img.tgz>

It constructs the edge order following Section 5 and discovers a merging opportunity following Section 7. We use Z3 [18] as the back end SMT and Max-SAT solver. As an additional contribution, our tool can also generate instances in the extended DIMACS format [3] allowing us to use any Max-SAT solvers. This output may also provide a source of industrially inspired examples for Max-SAT competition benchmarks.

Our benchmarks comprise 72 CSS files coming from three different sources (see Appendix D.2 for a complete listing):

- We collected CSS files from each of the top 20 websites on a global ranking list [1].
- CSS files were taken from a further 12 websites selected from the top 20-65 websites from the listing above.
- Finally, CSS files were taken from 11 smaller websites such as DBLP. These were examples used during the development of the tool.

This selection gives us a wide range of different types of websites, including large scale industrial websites and those developed by smaller teams. Note, several websites contained more than one CSS file and in this case we took a selection of CSS files from the site. Hence, we collected more examples than the number of websites used. Figure 13 gives the file-size distribution of the CSS files we collected.

In the following, we describe the optimisations implemented during the development of SATCSS. We then describe the particulars of the experimental setup and provide the results in Figure 15.

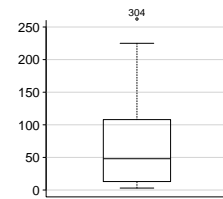


Figure 13: Box plot of the file sizes in kilobytes



## 8.1 Optimisations

We give an overview of the optimisations we used when implementing SATCSS. In section 8.4 we provide a detailed evaluation of the proposed optimisations.

- We introduce variables  $\bar{x}_k$  only for nodes appearing in the edge order, rather than for all nodes in a biclique. This is enough to be able to define sub-bicliques of any  $M_i$  that can appear at any position in the covering, but means that the smallest rule-merging opportunity cannot always be constructed. However, this reduces the search space and leads to an improvement to run times. For example, given a biclique

```
.a { color: red; background: blue }
```

where the property `color: red` appears in some pair of the edge ordering, but the property `background: blue` does not, we introduce a variable  $\bar{x}_k$  which is true whenever `color: red` is excluded from the biclique, but do not introduce a similar variable for `background: blue`. Since `background: blue` does not appear in the edge order, its presence can never cause a violation of the edge ordering. This is not the case for `color: red`, hence we still need to allow the possibility of removing it to satisfy the edge order.

- For a rule-merging application to reduce the size of the file, it must remove at least two nodes from the covering. Hence, for each  $M_i$  let  $j_l^i$  be the index of the *second* rule in  $\mathcal{C}$  containing some edge in  $M_i$  (not necessarily the same edge), and  $j_h^i$  be the index of the *last* rule containing some edge in  $M_i$ . Without loss of generality we assert

$$\bar{i}_M = i \Rightarrow (\bar{j} \geq j_l^i \wedge \bar{j} \leq j_h^i) .$$

- We performed the following optimisation when calculating  $(\{M_i\}_{i=1}^\mu, \mathcal{F})$ . The majority (58.4%) of benchmarks all maximal bicliques are orderable at all positions, and the mean percentage of maximal bicliques that were unorderable at some position was 0.39%, with a maximum of 5.84%. However, there are one or two of the largest examples of our experiments where this enumeration took a minute or so. Since the number of maximal bicliques that are unorderable at some position is so small, we decided in our implementation to simply remove all such bicliques from the analysis. In this case  $\{M_i\}_{i=1}^\mu$  is an enumeration of all maximal bicliques that are orderable at all positions  $j$ , and  $\mathcal{F}$  maps all values of  $[1, m]$  to  $\emptyset$ . This means that we may not be able to find the optimal rule-merging opportunity as not all bicliques are available, but since the amount of time required to find a rule-merging opportunity is reduced, we are able to find more merging opportunities to apply.
- Finally, we allowed a multi-threaded partitioned search. This works as follows. The search space is divided across  $n$  threads, and each thread partitions its search space into  $m$  partitions. During iteration  $j$ , thread  $k$  allows only those  $M_i$  where  $i = k * m + (j \bmod m)$ . If the fastest thread finds a merging opportunity in  $t$  seconds, we wait up to  $0.1t$  seconds for further instances. We take the best merging opportunity of those that have completed. A thread reports “no merging opportunities found” only if none of its partitions contain a merging opportunity.

For the experiments we implemented a simple heuristic to determine the number of threads and partitions to use. We describe this heuristic here, but first note that better results could likely be obtained with systematic parameter tuning rather than our ad-hoc settings. The heuristic used was to count the number of nodes in the CSS file; that is, the total number of selectors and property declarations (this total includes repetitions of the same node – we do not identify repetitions of

the same node). The tool creates enough partitions to give up to 750 nodes per partition. If only two partitions are needed, only one thread is used. Otherwise SATCSS first creates new threads – up to the total number of CPUs on the machine. Once this limit is reached, each thread is partitioned further until the following holds:  $number\ of\ threads * partitions\ per\ threads * 750 \leq number\ of\ nodes$ .

For the edge ordering:

- We do not support attribute selectors in full, but perform simple satisfiability checks on the most commonly occurring types of constraints. These cover all constraints in our benchmarks. However, we do support shorthand properties.
- Instead of doing a full Existential Presburger encoding, we do a backwards emptiness check of the CSS automata. This backwards search collects smaller Existential Presburger constraints describing the relationship between siblings in the tree, and checks they are satisfiable before moving to a parent node. Global constraints such as ID constraints are also collected and checked at the root of the tree. This algorithm is described in Appendix D.1.

## 8.2 Results

The experiments were run on a Dell Latitude e6320 laptop with 4Gb of RAM and four 2.7GHz Intel i7-2620M cores. The Python interpreter used was PyPy 5.6 [54] and the backend version of Z3 was 4.5. Each experiment was run up to a timeout of 30 minutes. In the case where CSS files used media queries (which our techniques do not yet support), we used stripmq [29] with default arguments to remove the media queries from CSS files. Also, we removed whitespaces, comments, and invalid CSS from all the CSS files before they were processed by the minifiers and our tool.

We used a timeout of 30 minutes because minification only needs to be applied once before deployment of the website. We note that the tool finds many merging opportunities during this period and a minified stylesheet can be produced after each application of rule-merging. This means the user will be able to trade time against the amount of size reduction made. Moreover, applications of rule-merging tend to show a “long-tail” where the first applications found provide larger savings and the later applications show diminishing returns (see Figure 14). The returns are diminishing because our MaxSAT encoding always searches for the rule merging opportunity with the largest saving.

## 8.3 Main Results

Table 1 and Figure 15 summarise the results. We compared our tool with six popular minifiers in current usage [60]. There are two groups of results: the first set show the results when either our tool or one of the minifiers is used alone on each benchmark; the second show the results when the CSS files are run first through a minifier and then through our tool. This second batch of results shows a significant improvement over running single tools in isolation. Thus, our tool complements and improves existing minification algorithms and our techniques may be incorporated into a suite of minification techniques.

Table 1 shows the savings, after whitespaces and comments are removed, obtained by SATCSS and the six minifiers when they are used alone and together. The table presents the savings in seven percentile ranks that include the minimal (0th), the median (50th), and the maximal values (100th). The upper half of the table shows the savings obtained in bytes, while the lower half shows the savings as percentages of the original file sizes. The first seven

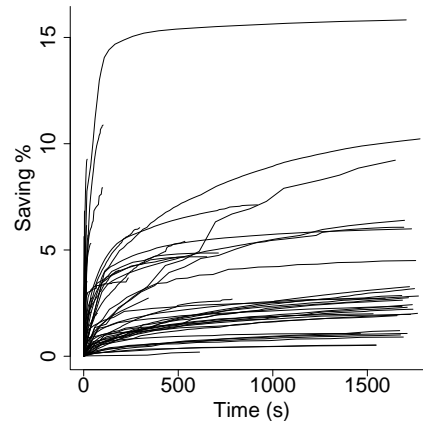


Figure 14: Savings against time for SATCSS on each benchmark

| Tool<br>Percentile | satcss | csso  | cssnano | cleancss | minify | yui  | cssmin | satcss +<br>csso | satcss +<br>cssnano | satcss +<br>cleancss | satcss +<br>minify | satcss +<br>yui | satcss +<br>cssmin |
|--------------------|--------|-------|---------|----------|--------|------|--------|------------------|---------------------|----------------------|--------------------|-----------------|--------------------|
| 0th                | 0      | -70   | 9       | 9        | 0      | 0    | 0      | 37               | 32                  | 45                   | 22                 | 24              | 24                 |
| 20th               | 320    | 258   | 100     | 199      | 5      | 13   | 3      | 467              | 444                 | 562                  | 362                | 366             | 369                |
| 40th               | 1011   | 877   | 462     | 738      | 28     | 35   | 24     | 1942             | 1353                | 1733                 | 1241               | 1082            | 1158               |
| 50th               | 2013   | 1527  | 1042    | 1354     | 72     | 92   | 61     | 4025             | 3022                | 4253                 | 2518               | 2290            | 2736               |
| 60th               | 3065   | 2955  | 1556    | 2081     | 153    | 157  | 117    | 5007             | 4370                | 4677                 | 3264               | 3279            | 3304               |
| 80th               | 4820   | 5656  | 3822    | 5449     | 355    | 326  | 212    | 8268             | 7202                | 8595                 | 5108               | 4765            | 5127               |
| 100th              | 81693  | 86367 | 76238   | 69573    | 3501   | 3464 | 3224   | 100734           | 97710               | 90710                | 78245              | 82101           | 77827              |
| 0th                | 0.00   | -0.56 | 0.02    | 0.10     | 0.00   | 0.00 | 0.00   | 0.14             | 0.12                | 0.33                 | 0.09               | 0.09            | 0.09               |
| 20th               | 1.79   | 1.45  | 0.75    | 1.21     | 0.02   | 0.04 | 0.01   | 2.96             | 2.65                | 3.60                 | 2.00               | 2.04            | 2.04               |
| 40th               | 2.90   | 2.54  | 1.47    | 2.55     | 0.10   | 0.17 | 0.09   | 4.90             | 3.94                | 5.05                 | 3.47               | 3.41            | 3.41               |
| 50th               | 3.79   | 3.29  | 1.84    | 3.03     | 0.15   | 0.21 | 0.18   | 5.78             | 4.90                | 6.10                 | 4.03               | 4.10            | 4.15               |
| 60th               | 4.62   | 4.40  | 2.42    | 3.64     | 0.22   | 0.26 | 0.21   | 6.50             | 5.64                | 7.45                 | 4.68               | 4.70            | 4.68               |
| 80th               | 8.12   | 6.13  | 5.18    | 6.04     | 0.56   | 0.64 | 0.49   | 10.11            | 10.98               | 10.42                | 8.21               | 8.09            | 8.39               |
| 100th              | 26.44  | 27.95 | 24.67   | 25.68    | 6.44   | 6.38 | 3.56   | 32.60            | 31.62               | 29.54                | 25.59              | 26.57           | 25.66              |

Table 1: Percentile ranks of the savings in bytes (above) and in percentages (below)

columns in the table show the savings when either our tool or one of the minifiers is used alone; the rest of the columns show the results when the CSS files are processed first by a minifier and then by our tool. Figure 15 shows the same data in visual form. It can be seen that SATCSS tends to achieve greater savings than each of the six minifiers on our benchmarks. Furthermore, even greater savings can be obtained when SATCSS is used in conjunction with any of the six minification tools. More precisely, when run individually, SATCSS achieves savings with a third quartile of 6.90% and a median value of 3.79%, while the six minifiers achieve savings with third quartiles and medians up to 5.45% and 3.29%, respectively. When we run SATCSS after running any one of these minifiers, the third quartile of the savings can be increased to 8.26% and the median to 4.70%. The additional gains obtained by SATCSS on top of the six minifiers (as a percentage of the original file size) have a third quartile of 5.03% and a median value of 2.80%. Moreover, the ratios of the percentage of savings made by SATCSS to the percentage of savings made by the six minifiers have third quartiles of at least 136% and medians of at least 48%. These figures clearly indicate a substantial proportion of extra space savings made by SATCSS. We comment in the next section on how our work may be integrated with existing tools.

In Figure 14 we plot for each benchmark, the savings made as time progresses. Each line represents one benchmark. Our algorithm repeatedly searches for and applies merging opportunities. Once one opportunity has been found, the search begins again for another. Hence, the longer SATCSS is run, the more opportunities can be found, and the more space saved. Since we search for the optimal merging opportunities first, we can observe a long-tail, with the majority of the savings being made early in the run. Hence, SatCSS could be stopped early while still obtaining a majority of the benefits. We further note that in 43 cases, SATCSS terminated once no more merging opportunities could be found. In the remaining 29 cases, the timeout was reached. If the timeout were extended, further savings may be found.

Finally, we remark on the validation of our tool. First, our model is built upon formal principles using techniques that are proven to maintain the CSS semantics. Moreover, our tool verifies that the output CSS is semantically equivalent to the input CSS. Thus we are confident that our techniques are also correct in practice. A reliable method for truly comparing whether the rendering of webpages using the original CSS and the minified CSS is identical is a matter for future work, for which the recent Visual Logic of Panckekha *et al.* [49] may prove useful. In lieu of a systematic validation, for each of our experiment inputs, we have visually verified that the rendering remains unchanged by the minification. Such a visual inspection can, of course, only be considered a sanity-check.

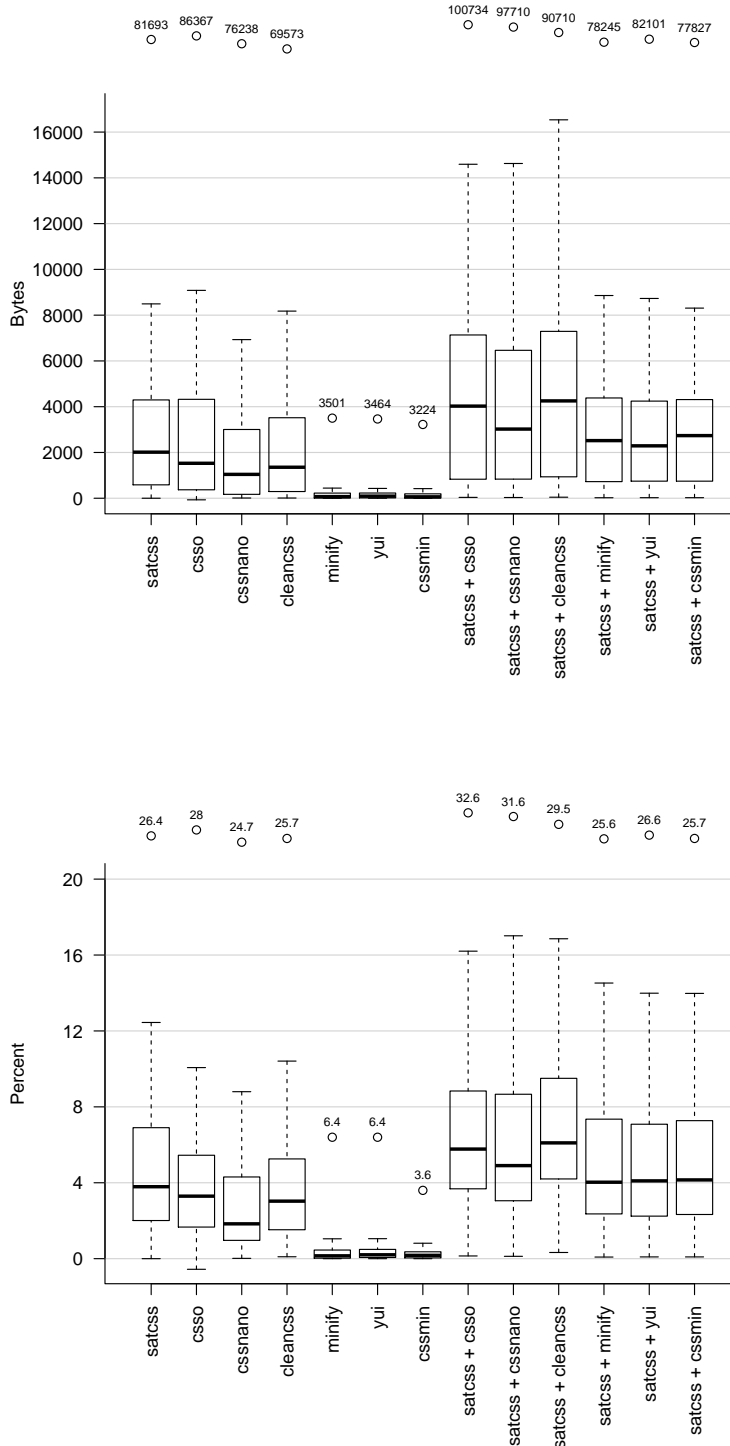


Figure 15: Box plots of the savings in bytes (above) and in percentages (below)

## 8.4 Evaluations of Optimisations

For each of our optimisations, we ran SATCSS on the unminified 72 benchmarks with the optimisation disabled. The results are shown in Figure 16 and Figure 17. In Figure 17, the “satcss” column shows the performance of SATCSS implemented as described above. As described above, of the 72 benchmarks, SATCSS completed within the timeout (no more merging-opportunities could be found) in 43 cases, and was stopped early due to the timeout in the remaining 29. The comparison with the effect of disabling optimisations is shown in Table 2.

Figure 16 shows the time taken by SATCSS to construct the edge order for each benchmark using the optimised emptiness of intersection algorithm presented in Appendix D.1 and the non-optimised encoding presented in Section 6. The timeout was kept at 30 minutes. The optimised algorithm completed the edge order construction on all benchmarks, while the timeout was reached in 19 cases by the unoptimised algorithm. A clear advantage can be seen for the optimised approach.

Next we compared the straightforward encoding of the rule-merging problem into Max-SAT discussed at the start of Section 7 with the biclique encoding which was the main topic of Section 7. SATCSS using the straightforward encoding appears as “nobiclques” in Figure 17. With the straightforward encoding, the tool completed within the timeout in 45 cases, and reached the timeout in the remaining 27. It can be seen that there is a clear benefit to the biclique encoding.

We also consider the optimisation that introduces variables  $\bar{x}_k$  only for nodes appearing in the edge order, rather than for all nodes in a biclique. Disabling this optimisation appears as “fullexc” in Figure 17. With this optimisation disabled, the tool completed within the timeout in 36 cases, and reached the timeout in the remaining 36. We can see a modest gain from this simple optimisation.

We then study the effect of removing all unorderable bicliques. The “unlimbicques” column in Figure 17 shows the effect of allowing  $(\{M_i\}_{i=1}^{\mu}, \mathcal{F})$  to be calculated in full. That is, unorderable bicliques are split into orderable sub-bicliques. With full biclique enumeration, the tool completed within the timeout in 42 cases, and reached the timeout in the remaining 30. This had only a small effect on performance, which is expected. The purpose of this limiting biclique generation in SATCSS is to prevent the rare examples of unorderable bicliques from having a large effect on performance in some cases.

The next optimisation in Figure 17, “nothread”, shows the performance of SATCSS with multi-threading and partitioning disabled. That is, the search space is not split up across several Max-SAT encodings. Without multi-threading or partitioning, the tool completed within the timeout in 40 cases, and reached the timeout in the remaining 32. There is a noticeable degradation in performance when this optimisation is disabled. The final column in Figure 17 studies the following hypothetical scenario. The reader may wonder whether CSS developers may be able to improve performance by specifying invariants of the documents to which the CSS will be applied. For example, the developer may specify that a node with class *a* will never also have class *b*. In this case pairs such as  $(.a, .b)$  can be removed from the edge-order. Thus, document-specific knowledge may improve performance by reducing the number of ordering constraints that need to be maintained. The column “noord”

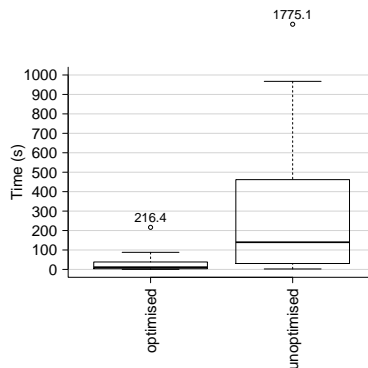


Figure 16: Box plots of the time taken to construct the CSS edge order with the optimised and unoptimised emptiness of intersection tests

| Optimisation | Terminated |
|--------------|------------|
| satcss       | 43         |
| nobiclques   | 45         |
| fullexc      | 36         |
| unlimbicques | 42         |
| nothread     | 40         |
| noord        | 50         |

Table 2: The number of examples terminating with no more discovered merging opportunities in standard mode and with certain optimisations disabled

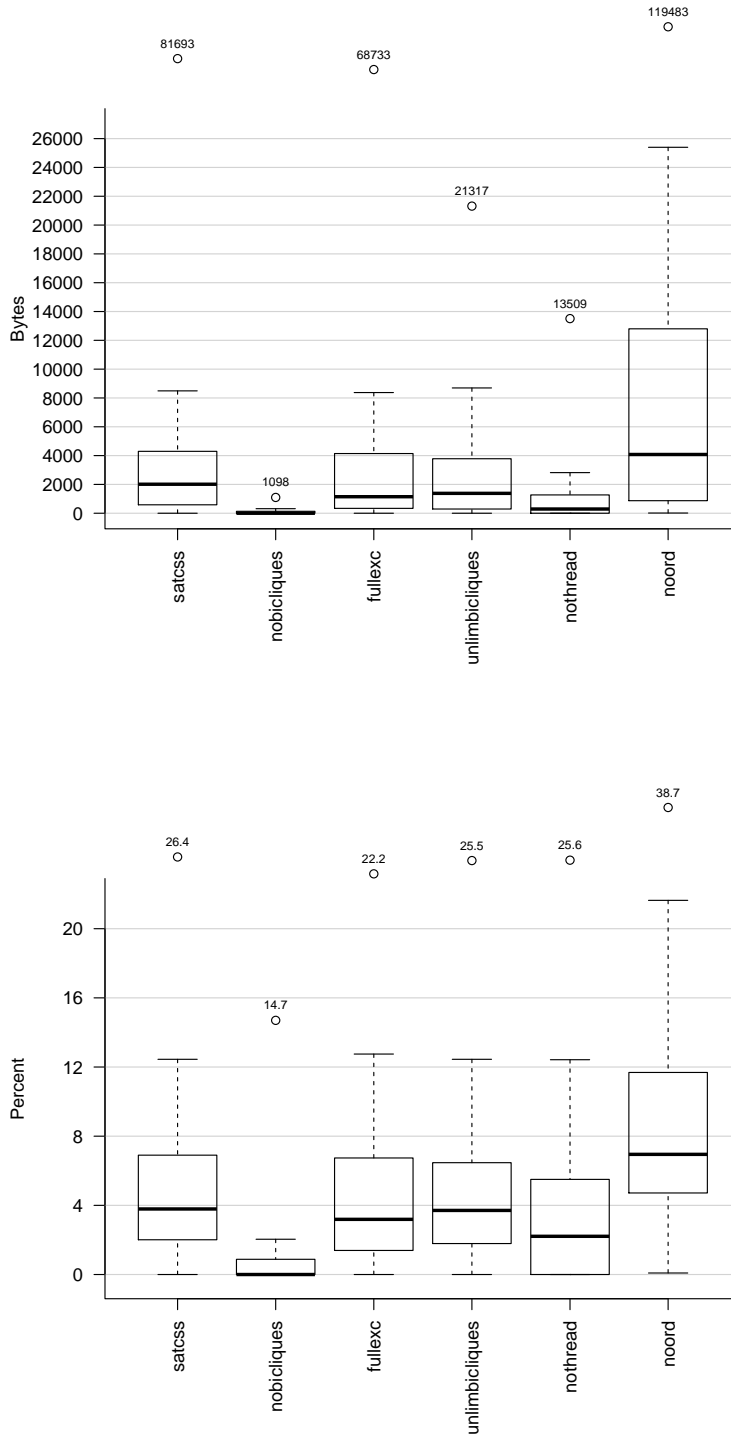


Figure 17: The savings in bytes (above) and in percentages (below) with certain optimisations disabled

shows the performance of SATCSS when the edge-ordering is empty, i.e., there is no edge ordering. Without the edge ordering, the tool completed within the timeout in 50 cases, and reached the timeout in the remaining 22. This represents a best-case scenario if document invariants were considered. Thus, if we were to extend SATCSS to also take a set of invariants, we could see a large improvement in the savings found. However, we also note that these savings do not dwarf the performance of SATCSS without invariants.

## 9 Related Work

CSS minification started to receive attention in the web programming community around the year 2000. To the best of our knowledge, the first major tools that could perform CSS minification were Yahoo! YUI Compressor [59] and Microsoft Ajax Minifier, both of which were developed around 2004–2006. This is followed by the development of many other CSS minifiers including (in no particular order) `cssmin` [8], `clean-css` [51], `csso` [19], `cssnano` [13], and `minify` [16]. Such minifiers mostly apply syntactic transformations including removing whitespace characters, comments, and replacing strings by their abbreviations (e.g. `#f60` by `#ff6600`). More and more advanced optimisations are also being developed. For example, `cssnano` provides a limited support of our rule-merging transformations, wherein only *adjacent* rules may be merged. The lack of techniques for handling the order dependencies of CSS rules [62] was most likely one main reason why a more general notion of rule-merging transformations (e.g. that can merge rules that are far away in the file) is not supported by CSS minifiers.

In our experiments, we ran SATCSS after running the existing minifiers described above. It is likely that the order of execution is important: the rewrites applied by existing minifiers will put the CSS into a more normalised format, which will improve the possibility of selectors sharing the same declarations. Moreover, these minifiers implement ad-hoc techniques, such as the limited rule-merging transformations of `cssnano` described above. After an application of our tool it is possible that some of these ad-hoc techniques may become applicable, leading to further savings. Thus, we posit that our techniques could be combined with the techniques of existing minifiers in a combined minification loop, which is run until a fixed point or timeout is reached.

Although the importance of CSS minification is understood in industry, the problem received little attention in academia until very recently. Below we will mention the small number of existing work on formalisation of CSS selectors and CSS minification, and other relevant work.

The lack of theories for reasoning about CSS selectors was first mentioned in the paper [25], wherein the authors developed a tree logic for algorithmically reasoning about CSS selectors, i.e., by developing algorithms for deciding satisfiability of formulas in the logic. This formalisation does *not* capture the whole class of CSS3 selectors; as remarked in their follow-up paper [11], the logic captures 70% of the selectors in their benchmarks from real-world CSS files. In particular, they do not fully support attribute selectors (e.g. `[l*="bob"]`). Our paper provides a *full* formalisation of CSS3 selectors. In addition, their tree logic can express properties that are *not* expressible in CSS. Upon a closer inspection, their logic is at least as expressive as  $\mu$ -calculus, which was a well-known logic in database theory for formalising query languages over XML documents, e.g., see [37, 47] for two wonderful surveys. As such, the complexity of satisfiability for their tree logic is EXPTIME-hard. Our formalisation captures *no more* than the expressive power of CSS3 selectors, which helps us obtain the much lower complexity NP and enables the use of highly-optimised SMT-solvers. There is a plethora of other work on logics and automata on unranked trees (with and without data), e.g., see [40, 6, 27, 67, 68, 66, 41, 31, 5, 48, 23, 38, 24, 21, 57, 17, 26] and the surveys [37, 47, 9]. However, none of these formalisms can capture certain aspects of CSS selectors (e.g. string constraints on attribute values), even though they are much more powerful than CSS selectors in other aspects (e.g. navigational).

There are a handful of research results on CSS minification that appeared in recent years. Loosely

speaking, these optimisations can be categorised into two: (a) *document-independent*, and (b) *document-dependent*. Document-independent optimisations are program transformations that are performed completely independently of the web documents (XML, HTML, etc.). On the other hand, document-dependent optimisations are performed with respect to a (possibly infinite) set of documents. Existing CSS minifiers *only* perform document-independent optimisations since they are meant to preserve the semantics of the CSS file *regardless* of the DOMs to which the CSS file is applied. Our work in this paper falls within this category too. Such optimisations are often the most sensible option *in practice* including (1) the case of *generic* stylesheets as part of web templates (e.g. WordPress), and (2) the case when the DOMs are generated by a program. Case (2) requires some explanation. A typical case of DOMs being generated by programs occurs in HTML5 web pages. An HTML5 application comes with a finite set of HTML documents, JavaScript code, and CSS files. The presence of JavaScript means that potentially *infinitely* many possible DOM-trees could be generated and displayed by the browser. Therefore, a CSS minification should *not* affect the rendering of *any* such tree by the browser. Although a document-dependent optimisation (that take these infinitely many trees into account) seems appropriate, this is far from realistic given the long-standing difficulty of performing sound static analysis for JavaScript especially in the presence of DOM-trees [56, 63, 2, 34, 33, 28]. This would make an interesting long-term research direction with many more advances on static analysis for JavaScript. However, the problem is further compounded by the multitude of frameworks deployed on the server-side for HTML creation (e.g. PHP, Java Server Pages, etc.), for which individual tools will need to be developed. Until then, a practical minifier for HTML5 applications will have to make do with document-independent optimisations for CSS files.

The authors of [43] developed a document-dependent dynamic analysis technique for detecting and removing unused CSS selectors in a CSS file that is part of an HTML5 application. A similar tool, called UnCSS [39], was also later developed. This is done by instrumenting the HTML5 application and removing CSS selectors that have not been used by the end of the instrumentation. The drawback of this technique is that it cannot test all possible behaviours of an HTML5 application and may accidentally delete selectors that can in reality be used by the application. It was noted in [28] that such tools may accidentally delete selectors, wherein the HTML5 application has event listeners that require user interactions. The same paper [28] develops a static analysis technique for overapproximating the set of generated DOM-trees by using tree rewriting for abstracting the dynamics of an HTML5 application. The technique, however, covers only a very small subset of JavaScript, and is difficult to extend without first overcoming the hard problem of static analysis of JavaScript.

The authors of [11] applied their earlier techniques [25] to develop a document-independent CSS minification technique that removes “redundant” property declarations, and merges two rules with semantically equivalent selectors. The optimisations that they considered are orthogonal to and can be used in conjunction with the optimisation that we consider in this paper. More precisely, they developed an algorithm for checking selector subsumption (given two selectors  $S_1$  and  $S_2$ , whether the selector  $S_1$  is subsumed by the selector  $S_2$ , written  $S_1 \subseteq S_2$ ). A redundant property declaration  $p$  in a rule  $R_1$  with a selector  $S_1$  can, then, be detected by finding a rule  $R_2$  that also contains the declaration  $p$  and has a selector  $S_2$  with a higher specificity than  $S_1$  and that  $S_1 \subseteq S_2$ . As another example, whenever we can show that the selectors  $S_1$  and  $S_2$  of two rules  $R_1$  and  $R_2$  to be semantically equivalent (i.e.  $S_1 \subseteq S_2$  and  $S_2 \subseteq S_1$ ), we may merge  $R_1$  with  $R_2$  under certain conditions. The authors of [11] provided sufficient conditions for performing this merge by relating the specificities of  $S_1$  and  $S_2$  with the specificities of other related selectors in the file (but not accounting for the order of appearances of these rules in the file). In general, a CSS rule might have multiple selectors (a.k.a. selector group), each with a different specificity, and it is not clear from the presentation of the paper [11] how their optimisations extend to the general case.



The authors of [42] developed a document-dependent<sup>10</sup> CSS minification method with an advanced type of rule merging as one of their optimisations. This is an ambitious work utilising a number of techniques from areas such as data mining and constraint satisfaction. Although their work differs from ours because of its document-dependence, the use of rule-merging is closely related to our own, hence we will describe in detail some key differences with our approach. The techniques presented in this paper can be viewed as a substantial generalisation of their rule merging optimisation. Loosely speaking, in terms of our graph-theoretic framework, their technique first enumerates all maximal bicliques with at least two selectors. This is done with the help of an association rule mining algorithm (from data mining) with a set of property declarations viewed as an *itemset*. Second, for each such maximal biclique  $B$ , a value  $n$  is computed that reflects how much saving will be obtained if  $B$  could somehow be inserted into the file and every occurrence of each property declaration in  $B$  is erased from the rest of the CSS file. Note that  $n$  is independent of where  $B$  is inserted into the CSS file. Third, for each such maximal biclique  $B$  (ranked according to their values in a non-increasing order), a solver for the (finite-domain) constraint satisfaction problem is invoked to check whether  $B$  can be placed in the file (with every occurrence of each property declaration in  $B$  is erased from the rest of the CSS file) while preserving the order dependency. If this check fails, the solver will also be invoked to check if one can insert sub-bicliques of  $B$  (with a maximal set  $S$  of selectors with  $|S| \geq 2$ ) in the file. Possible positions in the file to place each selector of  $B$  are encoded as variables constrained by the edge order dependency that is *relativised* to the provided HTML documents. To test whether two edges should be ordered in this relativised edge order, the selectors are not subject to a full intersection test, but instead a *relativised intersection* test that checks whether there is some node in the given finite set of html documents that is matched by both selectors. Their techniques do not work when the HTML documents are not given, which we handle in this paper. Another major difference to our paper is that their algorithm sequentially goes through every maximal biclique  $B$  (ranked according to their values) and checks if it can be inserted into the file, which is computationally too prohibitive especially when the (unrelativised) edge order  $\prec$  is used. Our algorithm, instead, fully relegates the search of an appropriate  $B$  and the position in the file to place it to a highly-optimised Max-SAT solver, which scales to real-world CSS files. In addition, their type of rule merging is also more restricted than ours for two other reasons. First, the new rule inserted into the file has to contain a maximal set of selectors. This prohibits many rule-merging opportunities and in fact does not subsume the merging adjacent rule optimisation of `cssnano` [13] in general. For example, consider the CSS file

```
.class1 { color:blue }
.class2 { color:blue }
.class3 { color:red }
.class4 { color:blue }
```

Notice that we cannot group together the first, second, and fourth rules since this would change the colour of a node associated with the classes `class2` and `class3`, or with the classes `.class3` and `class4`. On the other hand, the first two rules can be merged resulting in the new file

```
.class1, .class2 { color:blue }
.class3 { color:red }
.class4 { color:blue }
```

However, this is not permitted by their merging rule since `.class1, .class2{color:blue}` does not contain a maximal set of selectors. Second, given a maximal biclique  $B$ , their merging operation erases *every* occurrence of the declarations of  $B$  everywhere else in the file. This further rules out certain rule-merging opportunities. For example, consider the CSS file

<sup>10</sup>More precisely, dependent on a given finite set of HTML documents

```
.class1 { color:blue; font-size: large }
.class2 { color:blue; font-size: large }
.class4 { font-size: large }
.class3 { color:red }
.class4 { color:blue }
```

and observe the following maximal biclique in the file.

```
.class1, .class2, .class4 { color:blue; font-size: large }
```

Unfortunately, this is not a valid opportunity using their merging rule since this CSS file is not equivalent to

```
.class1, .class2, .class4 { color:blue; font-size: large }
.class3 { color:red }
```

nor to the following file.

```
.class3 { color:red }
.class1, .class2, .class4 { color:blue; font-size: large }
```

In this paper, we permit duplicate declarations, and would insert this maximal biclique just before the fourth rule in the file (and perform trim) resulting in the following equivalent file.

```
.class1, .class2, .class4 { color:blue; font-size: large }
.class3 { color:red }
.class4 { color: blue }
```

Finally, each maximal biclique in the enumeration of [42] does *not* allow two property declarations with the same property name. As we explained in Section 7, CSS rules satisfying this property are rather common since they support fallback options. Handling such bicliques (which we do in this paper) requires extra technicalities, e.g., the notion of orderable bicliques, and adding an order to the declarations in a biclique.

We also mention that there have been works [44, 50, 32] on solving web page layout using constraint solvers. These works are orthogonal to this paper. For example, [50] provides a mechanised formalisation of the semantics of CSS for web page layout (in quantifier-free linear arithmetic), which allows them to use an SMT-solver to automatically reason about layout. Our work provides a full formalisation of CSS selectors, which is not especially relevant for layout. Conversely, the layout semantics of various property declarations is not relevant in our CSS minification problem.

Finally, we also mention the potential application of rule-merging to CSS refactoring. This was already argued in [42], wherein the metric of minimal file size is equated with minimal redundancies. More research is required to discover other classes of CSS transformations and metrics that are more meaningful in the context of improving the design of stylesheets. Constraint-based refactoring has also been studied in the broader context of programming languages, e.g., see [64, 69]. It would be interesting to study how refactoring for CSS can be cast into the framework of constraint-based refactoring as previously studied (e.g. in [64]).

## 10 Conclusion and Future Work

We have presented a new CSS minification technique via merging similar rules. Our techniques can handle stylesheets composed of CSS rules which contain a set of CSS Level 3 selectors and list of property declarations. This technique exploits the fact that new rules may be introduced that render other

parts of the document redundant. After removing the redundant parts, the overall file size may be reduced. Such a solution has required the development of a complete formalisation of CSS selectors and their intersection problem as well as a formalisation of the dependency ordering present in a stylesheet. This intersection problem was solved by means of an efficient encoding to quantifier-free integer linear arithmetic, for which there are highly-optimised SMT solvers. Moreover, we have formalised our CSS rule-merging problem and presented a solution to this problem using an efficient encoding into MaxSAT formulas. These techniques have been implemented in our tool SATCSS which we have comprehensively compared with state-of-the-art minification tools. Our results show clear benefits of our approach.

Both our formalisation and our tool strictly follow the W3C specifications. In practice, web developers may not always follow these guidelines, and implement convenient abuses that do not trouble existing web browsers. One particular example is the use of ID values that are not necessarily unique. In this example case, it would be possible to treat ID values similarly to classes, and relax our analysis appropriately. In general, one may wish to adapt our constraints to handle other common abuses. However, this is beyond the scope of the current work.

CSS preprocessors such as Less [58] and Sass [14] — which extend the CSS language with useful features such as variables and partial rules — are commonly used in web development. Since Less and Sass code is compiled into CSS before deployment, our techniques are still applicable.

There are many technologies involved in website development and deployment. These technologies provide a variety of options for further research, some of which we briefly discuss here.

First, we may expand the scope of the CSS files we target. For example, we may expand our definition of CSS selectors to include features proposed in the CSS Selectors Level 4 working draft [20] (i.e. still not stable). These features include extensions of the negation operator to allow arbitrary selectors to be negated. It would be interesting to systematically investigate the impact of these features on the complexity of the intersection problem. We believe that such a systematic study would be informative in determining the future standards of CSS Selectors.

Another related technology is that of *media queries* that allow portions of a CSS file to only be applied if the host device has certain properties, such as a minimum screen size. This would involve defining semantics of media queries (not part of selectors), and extending our rule-merging problem to include media queries and rules grouped under media queries.

Second, we could also consider additional techniques for stylesheet optimisation. Currently we take a greedy approach where we search for the “best” merging opportunity at each iteration. Techniques such as simulated annealing allow a proportion of non-greedy steps to be applied (i.e. choose a merging opportunity that does not provide the largest reduction in file size). This allows the optimisation process to explore a larger search space, potentially leading to improved final results. Another approach might be to search for multiple simultaneous rule-merging opportunities.

Finally, our current optimisation metric is the raw file size. We could also attempt to provide an encoding that seeks to find the best file size reduction after gzip compression. [Gzip is now supported by many web hosts and most modern browsers (though not including old IE browsers).] One simple technique that could help bring down the compressed file size is to sort the selectors and declarations in each rule after the minification process is done [35].

## Acknowledgements

We are grateful for the support that we received from the Engineering and Physical Sciences Research Council [EP/K009907/1], Google (Faculty Research Award), and European Research Council (grant agreement no. 759969). We also thank Davood Mazinianian for answering questions about his work.

## References

- [1] Alexa Internet. Alexa top 500 global sites. <https://www.alexa.com/topsites>, 2017. Referred in April 2017.
- [2] E. Andreasen and A. Møller. Determinacy in static analysis for jQuery. In *OOPSLA*, pages 17–31, 2014.
- [3] Josep Argelich, Chu Min Li, Felip Manyà, and Jordi Planes. Max-sat’16 competition. <http://maxsat.ia.udl.cat/>, 2016. Referred in April 2017.
- [4] Tab Atkins Jr., Erika J. Etemad, and Florian Rivoal. Css snapshot 2017. <https://www.w3.org/TR/css-2017>, 2017. Referred in August 2017.
- [5] Michael Benedikt, Wenfei Fan, and Floris Geerts. XPath satisfiability in the presence of dtDs. In *Proceedings of the Twenty-fourth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 13-15, 2005, Baltimore, Maryland, USA*, pages 25–36, 2005.
- [6] Michael Benedikt, Wenfei Fan, and Gabriel M. Kuper. Structural properties of XPath fragments. In *Database Theory - ICDT 2003, 9th International Conference, Siena, Italy, January 8-10, 2003, Proceedings*, pages 79–95, 2003.
- [7] Nikolaj Bjørner and Nina Narodytska. Maximum satisfiability using cores and correction sets. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pages 246–252, 2015.
- [8] Johan Bleuzen. cssmin. <https://www.npmjs.com/package/cssmin>, 2017. Referred August 2017.
- [9] Mikołaj Bojańczyk. Automata for data words and data trees. In *Proceedings of the 21st International Conference on Rewriting Techniques and Applications, RTA 2010, July 11-13, 2010, Edinburgh, Scotland, UK*, pages 1–4, 2010.
- [10] Bert Bos. Cascading style sheets level 2 revision 2 (css 2.2) specification. <https://www.w3.org/TR/CSS22/>, 2016. Referred August 2017.
- [11] Martí Bosch, Pierre Genevès, and Nabil Layaïda. Reasoning with style. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pages 2227–2233, 2015.
- [12] Aaron R. Bradley and Zohar Manna. *The Calculus of Computation: Decision Procedures with Applications to Verification*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2007.
- [13] Ben Briggs and Contributors. cssnano. <http://cssnano.co>, 2015. Referred in August 2017.
- [14] Hampton Catlin, Natalie Weizenbaum, Chris Eppstein, and Contributors. Sass. <http://sass-lang.com/>, 2006. Referred in August 2017.
- [15] T. Çelik, E. J. Etemad, D. Glazman, I. Hickson, P. Linss, and J. Williams. Selectors level 3: W3c recommendation 29 september 2011. <http://www.w3.org/TR/2011/REC-css3-selectors-20110929/>, 2011. Referred in August 2017.
- [16] Steve Clay and Contributors. minify. <https://github.com/mrclay/minify>, 2017. Referred in August 2017.
- [17] Claire David, Leonid Libkin, and Tony Tan. Efficient reasoning about data trees via integer linear programming. *ACM Trans. Database Syst.*, 37(3):19:1–19:28, 2012.
- [18] L. Mendonça de Moura and N. Bjørner. Z3: An efficient SMT solver. In *TACAS*, 2008.
- [19] Roman Dvornov and Contributors. csso. <https://github.com/css/csso>, 2017. Referred in August 2017.
- [20] Erika J. Etemad and Tab Atkins Jr. Selectors level 4: W3c working draft 2 may 2013. <http://www.w3.org/TR/2013/WD-selectors4-20130502/>, 2013.
- [21] Diego Figueira. *Reasoning on words and trees with data: On decidable automata on data words and data trees in relation to satisfiability of LTL and XPath*. PhD thesis, Ecole Normale Supérieure de Cachan, 2010.
- [22] Ferenc Gécseg and Magnus Steinby. Handbook of formal languages, vol. 3. chapter Tree Languages, pages 1–68. Springer-Verlag New York, Inc., New York, NY, USA, 1997.
- [23] Floris Geerts and Wenfei Fan. Satisfiability of XPath queries with sibling axes. In *Database Programming*

- Languages, 10th International Symposium, DBPL 2005, Trondheim, Norway, August 28-29, 2005, Revised Selected Papers*, pages 122–137, 2005.
- [24] Pierre Genevès and Nabil Layaïda. A system for the static analysis of XPath. *ACM Trans. Inf. Syst.*, 24(4):475–502, 2006.
- [25] Pierre Genevès, Nabil Layaïda, and Vincent Quint. On the analysis of cascading style sheets. In *Proceedings of the 21st World Wide Web Conference 2012, WWW 2012, Lyon, France, April 16-20, 2012*, pages 809–818, 2012.
- [26] Pierre Genevès, Nabil Layaïda, Alan Schmitt, and Nils Gesbert. Efficiently deciding  $\mu$ -calculus with converse over finite trees. *ACM Trans. Comput. Log.*, 16(2):16:1–16:41, 2015.
- [27] Georg Gottlob and Christoph Koch. Monadic queries over tree-structured data. In *17th IEEE Symposium on Logic in Computer Science (LICS 2002), 22-25 July 2002, Copenhagen, Denmark, Proceedings*, pages 189–202, 2002.
- [28] Matthew Hague, Anthony Widjaja Lin, and C.-H. Luke Ong. Detecting redundant CSS rules in HTML5 applications: a tree rewriting approach. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2015, part of SPLASH 2015, Pittsburgh, PA, USA, October 25-30, 2015*, pages 1–19, 2015.
- [29] Jake Hickenlooper. stripmq. <https://www.npmjs.com/package/stripmq>, 2014. Referred in August 2017.
- [30] Ian Hickson, Robin Berjon, Steve Faulkner, Travis Leithead, Erika Doyle Navara, Edward O’Connor, and Silvia Pfeiffer. Html5. <https://www.w3.org/TR/html5/>, 2014. Referred August 2017.
- [31] Jan Hidders. Satisfiability of XPath expressions. In *Database Programming Languages, 9th International Workshop, DBPL 2003, Potsdam, Germany, September 6-8, 2003, Revised Papers*, pages 21–36, 2003.
- [32] Thibaud Hottelier and Rastislav Bodik. Synthesis of layout engines from relational constraints. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2015*, pages 74–88, New York, NY, USA, 2015. ACM.
- [33] S. H. Jensen, M. Madsen, and A. Møller. Modeling the HTML DOM and browser API in static analysis of JavaScript web applications. In *SIGSOFT/FSE*, pages 59–69, 2011.
- [34] S. H. Jensen, A. Møller, and P. Thiemann. Type Analysis for JavaScript. In *SAS*, pages 238–255, 2009.
- [35] Meitar Moscovitz Joseph R. Lewis. *AdvancED CSS*. APress, 2010.
- [36] Enver Kayaaslan. On enumerating all maximal bicliques of bipartite graphs. In *9th Cologne-Twente Workshop on Graphs and Combinatorial Optimization, Cologne, Germany, May 25-27, 2010. Extended Abstracts*, pages 105–108, 2010.
- [37] Leonid Libkin. Logics for unranked trees: An overview. *Logical Methods in Computer Science*, 2(3), 2006.
- [38] Leonid Libkin and Cristina Sirangelo. Reasoning about XML with temporal logics and automata. *J. Applied Logic*, 8(2):210–232, 2010.
- [39] Giacomo Martino and Contributors. Uncss. <https://github.com/giakki/uncss>, 2013. Referred April 2017.
- [40] Maarten Marx. Conditional xpath. *ACM Trans. Database Syst.*, 30(4):929–959, 2005.
- [41] Maarten Marx and Maarten de Rijke. Semantic characterizations of navigational xpath. *SIGMOD Record*, 34(2):41–46, 2005.
- [42] Davood Mazinanian, Nikolaos Tsantalis, and Ali Mesbah. Discovering refactoring opportunities in cascading style sheets. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, (FSE-22), Hong Kong, China, November 16 - 22, 2014*, pages 496–506, 2014.
- [43] Ali Mesbah and Shabnam Mirshokraie. Automated analysis of CSS rules to support style maintenance. In *34th International Conference on Software Engineering, ICSE 2012, June 2-9, 2012, Zurich, Switzerland*, pages 408–418, 2012.
- [44] Leo A. Meyerovich and Rastislav Bodík. Fast and parallel webpage layout. In *Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26-30, 2010*, pages 711–720, 2010.

- [45] A. Muscholl and I. Walukiewicz. An np-complete fragment of LTL. *Int. J. Found. Comput. Sci.*, 16(4):743–753, 2005.
- [46] Nina Narodytska and Fahiem Bacchus. Maximum satisfiability using core-guided maxsat resolution. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, July 27 -31, 2014, Québec City, Québec, Canada.*, pages 2717–2723, 2014.
- [47] Frank Neven. Automata theory for xml researchers. *SIGMOD Rec.*, 31(3):39–46, September 2002.
- [48] Frank Neven and Thomas Schwentick. On the complexity of xpath containment in the presence of disjunction, dtDs, and variables. *Logical Methods in Computer Science*, 2(3), 2006.
- [49] Pavel Panchekha, Adam T. Geller, Michael D. Ernst, Zachary Tatlock, and Shoaib Kamil. Verifying that web pages have accessible layout. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2018, Philadelphia, PA, USA, June 18-22, 2018*, pages 1–14, 2018.
- [50] Pavel Panchekha and Emina Torlak. Automated reasoning for web page layout. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2016, part of SPLASH 2016, Amsterdam, The Netherlands, October 30 - November 4, 2016*, pages 181–194, 2016.
- [51] Jakub Pawlowicz. clean-css. <https://github.com/jakubpawlowicz/clean-css>, 2017. Referred in August 2017.
- [52] René Peeters. The maximum edge biclique problem is np-complete. *Discrete Applied Mathematics*, 131(3):651–654, 2003.
- [53] Justyna Petke. *Bridging Constraint Satisfaction and Boolean Satisfiability*. Artificial Intelligence: Foundations, Theory, and Algorithms. Springer, 2015.
- [54] Armin Rigo and Contributors. Pypy. <http://pypy.org/>, 2007. Referred in August 2017.
- [55] B. Scarpellini. Complexity of subcases of presburger arithmetic. *Trans. of AMS*, 284(1):203–218, 1984.
- [56] M. Schäfer, M. Sridharan, J. Dolby, and F. Tip. Dynamic determinacy analysis. In *PLDI*, pages 165–174, 2013.
- [57] Helmut Seidl, Thomas Schwentick, Anca Muscholl, and Peter Habermehl. Counting in trees for free. In *Automata, Languages and Programming: 31st International Colloquium, ICALP 2004, Turku, Finland, July 12-16, 2004. Proceedings*, pages 1136–1149, 2004.
- [58] Alexis Sellier and Contributors. Less. <http://lesscss.org/>, 2009. Referred in August 2017.
- [59] Thomas Sha and Contributors. Yui compressor. <http://yui.github.io/yuicompressor/>, 2014. Referred August 2017.
- [60] Slant. Eight best css minifiers as of 2017. <https://www.slant.co/topics/261/~best-css-minifiers>, 2017. Referred in April 2017.
- [61] Jennifer Slegg. Google mobile first index: Page speed included as a ranking factor. *The SEM Post*, March 2017.
- [62] Steve Souders. *High Performance Web Sites: Essential Knowledge for Front-End Engineers*. O’Reilly Media, 2007.
- [63] M. Sridharan, J. Dolby, S. Chandra, M. Schäfer, and F. Tip. Correlation tracking for points-to analysis of JavaScript. In *ECOOP*, pages 435–458, 2012.
- [64] Friedrich Steimann. Constraint-based refactoring. *ACM Trans. Program. Lang. Syst.*, 40(1):2:1–2:40, 2018.
- [65] Larry J. Stockmeyer and Albert R. Meyer. Word problems requiring exponential time: Preliminary report. In *Proceedings of the 5th Annual ACM Symposium on Theory of Computing, April 30 - May 2, 1973, Austin, Texas, USA*, pages 1–9, 1973.
- [66] Balder ten Cate, Tadeusz Litak, and Maarten Marx. Complete axiomatizations for XPath fragments. *J. Applied Logic*, 8(2):153–172, 2010.
- [67] Balder ten Cate and Maarten Marx. Navigational XPath: calculus and algebra. *SIGMOD Record*, 36(2):19–26, 2007.

- [68] Balder ten Cate and Maarten Marx. Axiomatizing the logical core of XPath 2.0. *Theory Comput. Syst.*, 44(4):561–589, 2009.
- [69] Frank Tip, Robert M. Fuhrer, Adam Kiezun, Michael D. Ernst, Ittai Balaban, and Bjorn De Sutter. Refactoring using type constraints. *ACM Trans. Program. Lang. Syst.*, 33(3):9:1–9:47, 2011.
- [70] Unicode, Inc. The unicode standard, version 9.0. <http://www.unicode.org/versions/Unicode9.0.0>, 2016. Referred in August 2017.
- [71] Moshe Y. Vardi. An automata-theoretic approach to linear temporal logic. In *Logics for Concurrency - Structure versus Automata (8th Banff Higher Order Workshop, August 27 - September 3, 1995, Proceedings)*, pages 238–266, 1995.
- [72] Mihalis Yannakakis. Node- and edge-deletion np-complete problems. In *Proceedings of the 10th Annual ACM Symposium on Theory of Computing, May 1-3, 1978, San Diego, California, USA*, pages 253–264, 1978.

# Appendix

## A Additional Material for Max-SAT Encoding

Recall, given a valid covering  $\mathcal{C} = \{\overline{B}_i\}_{i=1}^m$  of a CSS graph  $\mathcal{G}$ , we aim to find a rule  $\overline{B} = (X, \overline{Y})$  and a position  $j$  that minimises the weight of  $\mathcal{C}[\overline{B} \rightarrow j]_{\downarrow}$ .

In this section we give material omitted from Section 7. We begin with a definition of orderability that will be useful for the remainder of the section. Then we will discuss how to produce the pair  $(\{M_i\}_{i=1}^m, \mathcal{F})$ . Finally we will show that our encoding is correct.

### A.1 Orderable Biclques

To insert a biclique  $B = (X, Y)$  into the covering, we need to make sure the order of its edges respects the edge order. We can only order the edges by ordering the properties in the biclique. More precisely, if we insert the biclique at position  $j$ , we need all edges in  $B$  that do not appear later in the file (i.e. in  $\{\overline{B}_i\}_{i=j+1}^m$ ) to respect the edge order. This is because it is only the last occurrence of an edge that influences the semantics of the stylesheet. Thus, let

$$E_j^B = \{e \in B \mid \text{index}(e) \leq j\} .$$

The edge ordering implies a required ordering of  $E_j^B$ , which implies an ordering on the properties in  $Y$ . This ordering is defined as follows. For all  $p_1, p_2 \in Y$  we have

$$p_1 \ll_j^B p_2 \Leftrightarrow \exists (s_1, p_1), (s_2, p_2) \in E_j^B . (s_1, p_1) \prec^* (s_2, p_2) .$$

That is, we require  $p_1$  to appear before  $p_2$  if there are two edges  $(s_1, p_1)$  and  $(s_2, p_2)$  in  $B$  that must be ordered according to the transitive closure of  $\prec$ . A biclique is orderable iff its properties can be ordered in such a way to respect  $\ll_j^B$ .

**Definition A.1** (Orderable Biclques). *The biclique  $B$  is orderable at  $j$  if  $\ll_j^B$  is acyclic. That is, there does not exist a sequence  $(s_1, p_1), \dots, (s_n, p_n)$  such that  $(s_i, p_i) \ll_j^M (s_{i+1}, p_{i+1})$  for all  $1 \leq i < n$  and  $(s_1, p_1) = (s_n, p_n)$ .*

This can be easily checked in polynomial time. Moreover, if a biclique is orderable at a given position, a suitable ordering can be found by computing  $\ll_j^B$ , also in polynomial time. Thus, this proves Proposition 7.1.

## A.2 Enumerating Maximal Rules

Fix a covering  $\mathcal{C} = \{\overline{B}_i\}_{i=1}^m$  of a CSS graph  $\mathcal{G} = (S, P, E, \prec, \text{wt})$ . We show how to efficiently create the pair  $(\{M_i\}_{i=1}^m, \mathcal{F})$ . In fact, we will create the pair  $(\{M_i\}_{i=1}^m, \mathcal{F}^{1\text{st}})$  where  $\mathcal{F}^{1\text{st}}(j)$  is set of  $M_i$  for which  $j$  is the smallest position such that  $M_i$  is unorderable at position  $j$ . Computing  $\mathcal{F}$  from this is straightforward (though unnecessary since our encoding uses  $\mathcal{F}^{1\text{st}}$  directly).

Our algorithm begins with an enumeration of the maximal bicliques that are orderable at position 0, and iterates up to position  $m$ , extending the enumeration at each step and recording in  $\mathcal{F}^{1\text{st}}$  which maximal bicliques become unorderable at each position. In fact, we will construct a pair  $(\mathcal{M}, \mathcal{F}^{1\text{st}})$  where  $\mathcal{M}$  is a set of bicliques rather than a sequence. To construct a sequence we apply any ordering to the elements of  $\mathcal{M}$ .

### A.2.1 Initialisation

At position 0, all maximal bicliques  $(X, Y)$  with  $X \times Y \subseteq E$  are orderable. This is because all edges appearing in  $(X, Y)$  also appear in  $\{\overline{B}_i\}_{i=1}^m$ , hence, the semantics of  $\mathcal{C}[\overline{B} \rightarrow 0]$  will be decided by edges already in  $\mathcal{C}$ , and thus the ordering of the properties does not matter. Hence, we begin with a set of all maximal bicliques  $\mathcal{M}_0$ . That is, all bicliques  $(X, Y)$  such that  $X \times Y \subseteq E$  where there is no  $(X', Y')$  with  $X' \times Y' \subseteq E$  and  $X \times Y \subset X' \times Y'$ . Algorithms for generating such an enumeration are known. For example, we use the algorithm from Kayaaslan [36].

For an initial value  $\mathcal{F}_0^{1\text{st}}$  of  $\mathcal{F}^{1\text{st}}$ , we can simply take the empty function  $\emptyset$ .

### A.2.2 Iteration

Assume we have generated  $(\mathcal{M}_j, \mathcal{F}_j^{1\text{st}})$ . We show how to generate the extension  $(\mathcal{M}_{j+1}, \mathcal{F}_{j+1}^{1\text{st}})$ .

The idea is to find all elements of  $\mathcal{M}_j$  that are not orderable at position  $j + 1$ . Let  $\chi$  be this set. We first define

$$\mathcal{F}_{j+1}^{1\text{st}} = \mathcal{F}_j^{1\text{st}} \cup \{(j + 1, \chi)\} .$$

Then, for each biclique  $M \in \chi$ , we search for smaller bicliques contained within  $M$  that are maximal and orderable at position  $j + 1$ . This results in the extension  $\{M_i\}_{i=1}^{j+1}$  giving us  $(\{M_i\}_{i=1}^{j+1}, \mathcal{F}_{j+1}^{1\text{st}})$ .

Thus, the problem reduces to finding bicliques contained within some  $M$  that are maximal and orderable at position  $j + 1$ . We describe a simple algorithm for this in the next section.

### A.2.3 Algorithm for $(\mathcal{M}, \mathcal{F}^{1\text{st}})$

We write  $\text{Orderable}(M, j)$  to assert that a biclique  $M$  is orderable at position  $j$ . For now, assume we have the subroutine  $\text{OrderableSub}(M, j)$  which returns a set  $\mathcal{M}'$  of all orderable maximal bicliques at position  $j$  contained within  $M$ . We use the following algorithm to generate  $(\mathcal{M}, \mathcal{F}^{1\text{st}})$ .

```

 $\mathcal{M} := \mathcal{M}_0$ 
 $\mathcal{F}^{1\text{st}} := \emptyset$ 
for  $j := 1$  to  $m$  do
   $\chi := \emptyset$ 
  for all  $M \in \mathcal{M}$  do
    if  $\neg \text{Orderable}(M, j)$  then
       $\chi := \chi \cup \{M\}$ 
       $\mathcal{M} := \mathcal{M} \cup \text{OrderableSub}(M, j)$ 
    end if
  end for
   $\mathcal{F}^{1\text{st}} := \mathcal{F}^{1\text{st}} \cup \{(j + 1, \chi)\}$ 

```



**end for**  
**return**  $(\mathcal{M}, \mathcal{F}^{1st})$

Note, we can improve the algorithm by restricting the nested for all loop over elements  $M \in \mathcal{M}$  to only those  $M$  that are not orderable at  $m$ . This is because an ordering at position  $m$  is also an ordering at position  $j \leq m$ . Hence, these bicliques will never be unorderable and do not need to be checked repeatedly.

#### A.2.4 Generating Orderable Sub-Bicliques

We now give an algorithm for implementing the subroutine  $\text{OrderableSub}(M, j)$ . Naively we can simply generate all sub-bicliques  $M'$  of  $M$  and check  $\text{Orderable}(M', j)$ . However, to avoid the potentially high cost of such an iteration, we first determine which selectors and properties contribute to  $\ll_j^M$ . Removing nodes outside of these sets will not affect the orderability, hence we do not need to try removing them. Then we first attempt only removing one node from this set, computing all sub-bicliques that have one fewer element and are orderable. Then, for all nodes for which this fails, we attempt to remove two nodes, and so on. Note, if removing node  $w$  renders  $M$  orderable, we do not need to test any bicliques obtained by removing  $w$  and some other node  $w'$ , since this will not result in a maximal biclique.

Hence, we define the sets of candidate selectors and properties that may be removed to restore orderability. These are all selectors and nodes that contribute to  $\ll_j^M$ . That is

$$\Delta = \{s_1, s_2, p_1, p_2 \mid \exists (s_1, p_1), (s_2, p_2) \in E_j^M . (s_1, p_1) \prec^* (s_2, p_2)\} .$$

We define  $\text{OrderableSub}(M, j) = \text{OrderableSub}(M, j, \Delta)$  where  $\text{OrderableSub}(M, j, \Delta)$  generates a set  $\Omega$  of orderable sub-bicliques and is defined below. When  $M = (X, Y)$  we will abuse notation and write  $M \setminus \{w\}$  for  $(X \setminus \{w\}, Y)$  when  $w$  is a selector, and  $(X, Y \setminus \{w\})$  when  $w$  is a property. When defining the algorithm, we will collect all orderable bicliques in a set  $\Omega$ . We will further collect in  $\Delta'$  the set of all nodes which fail to create an orderable biclique when removed by themselves. We define  $\text{OrderableSub}(M, j, \Delta)$  recursively, where the recursive call attempts the removal of an increasing number of nodes. It is

```

 $\Omega := \emptyset$ 
 $\Delta' := \emptyset$ 
for all  $w \in \Delta$  do
   $M' := M \setminus \{w\}$ 
  if  $\text{Orderable}(M', j)$  then
     $\Omega := \Omega \cup \{M'\}$ 
  else
     $\Delta' := \Delta' \cup \{w\}$ 
  end if
end for
for all  $w \in \Delta'$  do
   $\Omega := \Omega \cup \text{OrderableSub}(M \setminus \{w\}, j, \Delta' \setminus \{w\})$ 
end for
return  $\Omega$ 

```

### A.3 Correctness of the Encoding

We argue Proposition 7.3 which claims that the encoding  $(\Pi_H, \Pi_S)$  is correct. To prove this we need to establish three facts.

1. If  $(\bar{B}, j)$  is a valid merging opportunity of  $\mathcal{C}$  and  $\bar{B} = (B, \triangleleft)$ , then  $(B, j)$  is a solution to the hard constraints.
2. If  $(\bar{B}, j)$  is generated from a solution to the hard constraints, it is a valid merging opportunity.
3. The weight of a solution generating  $(\bar{B}, j)$  is the size of  $\mathcal{C}[\bar{B} \rightarrow j]_{\downarrow}$ .

We argue these properties below.

1. Take a valid merging opportunity  $(\bar{B}, j)$  and let  $\bar{B} = (B, \triangleleft)$ . We construct a solution to the hard constraints. First, we assign  $\bar{j} = j$ . Next, since the merging opportunity is valid, we know  $B$  contains only edges in  $E$ . That is, it is contained within a maximal biclique. Furthermore, since  $\mathcal{C}[\bar{B} \rightarrow j]$  is valid, we know that  $B$  is orderable at position  $j$ . Thus, it is a sub biclique of some  $M_i$  in  $(\{M_i\}_{i=1}^{\mu}, \mathcal{F}^{1st})$ , and, moreover, it is not the case that  $M_i \in \mathcal{F}^{1st}(j')$  for some  $j' \leq j$ . Thus, we assign  $i_M = i$  and we know that

$$\bigwedge_{1 \leq j \leq m+1} \left( (\bar{j} \geq j) \Rightarrow \bigwedge_{M_i \in \mathcal{F}^{1st}(j)} (\bar{i}_M \neq i) \right)$$

is satisfied.

Additionally, for all  $w$  appearing in  $B$  but not in  $M_i$ , we set  $\bar{x}_{\rho_i(w)}$  to false, otherwise we set it to true. Thus  $\text{HasEdge}((s, p))$  holds only if  $(s, p)$  is an edge in  $B$ .

Next, we argue

$$\left( \bigwedge_{(s_1, p_1) \prec (s_2, p_2)} \text{HasEdge}((s_1, p_1)) \Rightarrow (\bar{j} \leq \text{index}((s_2, p_2)) \vee \text{HasEdge}((s_2, p_2))) \right)$$

is satisfied. This follows from  $\mathcal{C}[\bar{B} \rightarrow j]$  being valid. To see this, take some  $(s_1, p_1) \prec (s_2, p_2)$ . If  $(s_1, p_1)$  does not appear in  $B$ , then there is nothing to prove. If it does, we know  $(s_2, p_2)$  must appear later in the file. There are two cases. If  $j \leq \text{index}((s_2, p_2))$  then the clause is satisfied. Otherwise we must have  $(s_2, p_2)$  in  $B$  or edge order would be violated. Thus the clause also holds in this case.

2. We need to prove that if the hard constraints are satisfied, then then generated merging opportunity  $(\bar{B}, j)$  is valid. Let  $\bar{B} = (B, \triangleleft)$ . For  $\mathcal{C}[\bar{B} \rightarrow j]$  to be valid, we first have to show that  $B$  introduces no new edges to the stylesheet. This is immediate since  $B$  is a sub biclique of some  $M_i$  in  $(\{M_i\}_{i=1}^{\mu}, \mathcal{F}^{1st})$ , which can only contain edges in  $E$ .

Next, we need to argue that we can create the ordering  $\triangleleft$  for the properties in  $B$ . First note that  $M_i$  is orderable at position  $j$ . In particular, for any  $(s_1, p_1) \prec^* (s_2, p_2)$  with  $(s_1, p_1)$  and  $(s_2, p_2)$  appearing in  $M_i$ , we have  $p_1 \ll_j^{M_i} p_2$ . Since all edges in  $B$  also appear in  $M_i$ , the existence of an ordering is immediate.

Finally, we need to argue that  $\mathcal{C}[\bar{B} \rightarrow j]$  respects the edge order. Suppose  $(s_1, p_2) \prec (s_2, p_2)$ . To violate this ordering, we need to introduce a copy of  $(s_1, p_1)$  after the last copy of  $(s_2, p_2)$ . Thus, we must have  $(s_1, p_1)$  in  $B$ . However, from

$$\left( \bigwedge_{(s_1, p_1) \prec (s_2, p_2)} \text{HasEdge}((s_1, p_1)) \Rightarrow (\bar{j} \leq \text{index}((s_2, p_2)) \vee \text{HasEdge}((s_2, p_2))) \right)$$

we are left with two cases. In the first  $j \leq \text{index}((s_2, p_2))$  and the edge order is maintained. In the second, we also have  $(s_2, p_2)$  in  $B$ . However, the edge order is maintained because  $B$  is orderable. Thus we are done.

3. Finally, we argue that the weight of a satisfying assignment accurately reflects the size of  $\mathcal{C}[\overline{B} \rightarrow j]_{\downarrow}$ . This is fairly straightforward. The size of  $\mathcal{C}[\overline{B} \rightarrow j]_{\downarrow}$  comprises two parts: the size of  $\overline{B}$ , and the size of  $\mathcal{C}$  after the trim operation. It is immediate to see that the size of  $\overline{B}$  is equal to the size of all of its nodes. In particular, this is the size of all nodes of  $M_i$  that appear in  $\overline{B}$ . That is, have not been excluded. Thus the clause with weight  $\text{wt}(w)$

$$(\bar{i}_M = i) \Rightarrow \bar{x}_{\rho_i(w)} .$$

for each  $w$  appearing in  $M_i$  accurately computes the size of  $\overline{B}$ .

For the size of  $\mathcal{C}$  after the trim operation, we first use the assumption that  $\mathcal{C}$  has already been trimmed before applying the merging opportunity. Thus, any further nodes removed in  $\mathcal{C}[\overline{B} \rightarrow j]_{\downarrow}$  from a rule  $\overline{B}_{i'}$  must be removed because some edge  $e$  in  $\overline{B}$  also appears in  $\overline{B}_{i'}$  and, moreover, it was the case  $i' = \text{index}(e)$  and  $i' \leq j$ . In particular, we can only remove a node  $w$  from  $\overline{B}_{i'}$  if all edges  $e$  incident to  $w$  with  $i' = \text{index}(e)$  have  $e$  appearing in  $\overline{B}$  (else there will still be some edge preventing  $w$  from being trimmed after applying the merging opportunity). Thus, for each selector node  $s$ , we know it is not removed if the clause with weight  $\text{wt}(s)$

$$i \leq \bar{j} \wedge \bigwedge_{\substack{\text{index}((s,p))=i \\ p \in \overline{Y}}} \text{HasEdge}((s,p))$$

is not satisfied. Similarly for property nodes  $p$ . Thus, these clauses accurately count the size of the covering after trimming.

## B Additional Material for Section 5

### B.1 Handling Pseudo-Elements

CSS selectors can also finish with a *pseudo-element*. For example  $\varphi :: \text{before}$ . These match nodes that are not formally part of a document tree. In the case of  $\varphi :: \text{before}$  the selector matches a phantom node appearing before the node matched by  $\varphi$ . These can be used to insert content into the tree for stylistic purposes. For example

```
.a::before { content:">" }
```

places a “>” symbol before the rendering of any node with class  $a$ .

We divide CSS selectors into five different types depending on the pseudo-element appearing at the end of the selector. We are interested here in the nodes matched by a selector. The pseudo-elements  $:: \text{first-line}$ ,  $:: \text{first-letter}$ ,  $:: \text{before}$ , and  $:: \text{after}$  essentially match nodes inserted into the DOM tree. The CCS3 specification outlines how these nodes should be created. For our purposes we only need to know that the five syntactic cases in the above grammar can never match the same inserted node, and the selectors  $:: \text{first-letter}$  and  $:: \text{first-line}$  require that the node matched by  $\varphi$  is not empty.

Since we are interested here in the non-emptiness and non-emptiness-of-intersection problems, we will omit pseudo-elements in the remainder of this article, under the assumptions that

- selectors of the form  $\varphi :: \text{first-line}$  or  $\varphi :: \text{first-letter}$  are replaced by a selector  $\varphi : \text{not} (: \text{empty})$ , and
- selectors of the form  $\varphi, \varphi :: \text{before}$ , or  $\varphi :: \text{after}$  are replaced by  $\varphi$ , and
- we *never* take the intersection of two selectors  $\varphi$  and  $\varphi'$  such that it's not the case that either
  - $\varphi$  and  $\varphi'$  were derived from selectors containing no pseudo-elements, or
  - $\varphi$  and  $\varphi'$  were derived from selectors ending with the same pseudo-element.

In this way, we can test non-emptiness of a selector by testing its replacement. For non-emptiness-of-intersection, we know if two selectors end with different pseudo-elements (or one does not contain a pseudo-element, and one does), their intersection is necessarily empty. Thus, to check non-emptiness-of-intersection, we immediately return “empty” for any two selectors ending with different pseudo-elements. To check two selectors ending with the same pseudo-element, the problem reduces to testing the intersection of their replacements.

## B.2 NP-hardness of Theorem 5.2

**Lemma B.1.** *Given a CSS selector  $\varphi$ , deciding  $\exists T, \eta. T, \eta \models \varphi$  is NP-hard.*

*Proof.* We give a polynomial-time reduction from the NP-complete problem of non-universality of unions of arithmetic progressions [65, Proof of Theorem 6.1]. To define this, we first fix some notation. Given a pair  $(\alpha, \beta) \in \mathbb{N} \times \mathbb{N}$ , we define  $\llbracket (\alpha, \beta) \rrbracket$  to be the set of natural numbers of the form  $\alpha n + \beta$  for  $n \in \mathbb{N}$ . That is,  $\llbracket (\alpha, \beta) \rrbracket$  represents an arithmetic progression, where  $\alpha$  represents the *period* and  $\beta$  represents the *offset*. Let  $E \subseteq \mathbb{N} \times \mathbb{N}$  be a finite subset of pairs  $(\alpha, \beta)$ . We define  $\llbracket E \rrbracket = \bigcup_{(\alpha, \beta) \in E} \llbracket (\alpha, \beta) \rrbracket$ . The NP-complete problem is: given  $E$  (where numbers may be represented in unary or in binary representation), is  $\llbracket E \rrbracket \neq \mathbb{N}$ ? Observe that this problem is equivalent to checking whether  $\llbracket E+1 \rrbracket \neq \mathbb{N}_{>0}$  where  $E+1$  is defined by adding 1 to the offset  $\beta$  of each arithmetic progression  $(\alpha, \beta)$  in  $E$ . By complementation, this last problem is equivalent to checking whether  $\mathbb{N}_{>0} \setminus \llbracket E+1 \rrbracket \neq \emptyset$ . Since  $\mathbb{N}_{>0} \setminus \llbracket E+1 \rrbracket = \bigcap_{(\alpha, \beta) \in E} \llbracket \alpha, \beta+1 \rrbracket$ , the problem can be seen to be equivalent to testing the non-emptiness of

$$* \{ : \text{not} (: \text{nth-child}(\alpha n + (\beta + 1))) \mid (\alpha, \beta) \in E \} .$$

Thus, non-emptiness is NP-hard. □

## B.3 Handling !important and shorthand property names

Our approach handles the !important keyword and shorthand property names. In this section we explain the steps we take to account for them.

### B.3.1 The !important Keyword

First, the keyword !important in property declaration as is used in the rule

```
div { color:red !important }
```

can be used to override the cascading behaviour of CSS, e.g., in our example, if a node is matched by `div`, as well as a later rule  $R$  that assigns a different color, then assign `red` to `color` (unless  $R$  also has the keyword !important next to its color property declaration). To handle this, we can

extend the notion of specificity of a selector to the notion of specificity of a pair  $(s, p)$  of selector and property declaration, after which we may proceed as before (i.e. relating only two edges with the same specificity). Recall from [15] that the specificity of a selector is a 3-tuple  $(a, b, c) \in \mathbb{N} \times \mathbb{N} \times \mathbb{N}$  where  $a$ ,  $b$ , and  $c$  can be obtained by calculating the sum of the number of IDs, classes, tag names, etc. in the selector. Since the lexicographic order is used to rank the elements of  $S := \mathbb{N} \times \mathbb{N} \times \mathbb{N}$ , the specificity of a pair  $(s, p)$  can now be defined to be  $(i, a, b, c)$ , where  $(a, b, c)$  is the specificity of  $s$ , and  $i = 1$  if `!important` can be found in  $p$  (otherwise,  $i = 0$ ). In particular, this also handles the case where multiple occurrences of `!important` is found in the CSS file.

### B.3.2 Shorthand Property Names

*Shorthand property names* [10] can be used to simultaneously set the values of related property names. For example, `border: 5px solid red` is equivalent to

```
border-width: 5px; border-style: solid; border-color:red
```

In particular, this implies that  $(s, p)$  and  $(s, p')$  can be related in  $\prec$  if  $p$  defines `border`, while the other property  $p'$  defines `border-width`. One way to achieve this is to simply list all pairs of comparable property names, which can be done since only around 100 property names are currently officially related. [Incidentally, a close enough approximation is that one property name is a *prefix* of the other property name (e.g., `border` is a prefix of `border-style`), but this is not complete (e.g. `font` can be used to define `line-height`)]

## C Additional Material for Section 6

### C.1 Correctness of $\mathcal{A}_\varphi$ in Proposition 6.1

We show both soundness and completeness of  $\mathcal{A}_\varphi$ .

**Lemma C.1.** *For each CSS selector  $\varphi$  and tree  $T$ , we have*

$$(T, \eta) \in \mathcal{L}(\mathcal{A}_\varphi) \Rightarrow T, \eta \models \varphi.$$

*Proof.* Suppose  $(T, \eta) \in \mathcal{L}(\mathcal{A}_\varphi)$ . By construction of  $\mathcal{A}_\varphi$  we know that the accepting run must pass through all states  $\circ_1, \dots, \circ_n$  where  $\varphi = \sigma_1 \circ_1 \cdots \circ_{n-1} \sigma_n$ . Notice, in order to exit each state  $\circ_i$  a transition labelled by  $\sigma_i$  must be taken. Let  $\eta_i$  be the node read by this transition, which necessarily satisfies  $\sigma_i$ . Observe  $\eta_n = \eta$ . We proceed by induction. We have  $\eta_1$  satisfies  $\sigma_1$ . Hence, assume  $\eta_i$  satisfies  $\sigma_1 \circ_1 \cdots \circ_{i-1} \sigma_i$ . We show  $\eta_{i+1}$  satisfies  $\sigma_1 \circ_1 \cdots \circ_i \sigma_{i+1}$ .

We case split on  $\circ_i$ .

- When  $\circ_i = \gg$  we need to show  $\eta_{i+1}$  is a descendant of  $\eta_i$ . By construction of  $\mathcal{A}_\varphi$  the run reaches  $\eta_{i+1}$  in one of two ways. If it is via a single transition  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \circ_{i+1}$  then  $\eta_{i+1}$  is immediately a descendant of  $\eta_i$ . Otherwise the first transition is  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \bullet_i$ . The reached node is necessarily a descendant of  $\eta_i$ . To reach  $\eta_{i+1}$  a path is followed applying  $\rightarrow_+$  and  $\downarrow$  arbitrarily, which cannot reach a node that is not a descendant of  $\eta_i$ . Finally, the transition to  $\eta_{i+1}$  is via  $\rightarrow$  or  $\downarrow$  and hence  $\eta_{i+1}$  must also be a descendant of  $\eta_i$ .
- When  $\circ_i = >$  we need to show  $\eta_{i+1}$  is a descendant of  $\eta_i$ . By construction of  $\mathcal{A}_\varphi$  the run reaches  $\eta_{i+1}$  in one of two ways. If it is via a single transition  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \circ_{i+1}$  then  $\eta_{i+1}$  is immediately a

child of  $\eta_i$ . Otherwise the first transition is  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \bullet_i$ . The reached node is necessarily a child of  $\eta_i$ . To reach  $\eta_{i+1}$  only transitions labelled  $\rightarrow_+$  and  $\rightarrow$  can be followed. Hence, the node reached must also be a child of  $\eta_i$ .

- When  $o_i = +$  we need to show  $\eta_{i+1}$  is the next neighbour of  $\eta_i$ . Since the only path is a single transition labelled  $\rightarrow$  the result is immediate.
- When  $o_i = \sim$  we need to show  $\eta_{i+1}$  is a sibling of  $\eta_i$ . By construction of  $\mathcal{A}_\varphi$  the run reaches  $\eta_{i+1}$  in one of two ways. If it is via a single transition  $\circ_i \xrightarrow[\sigma_i]{\rightarrow} \circ_{i+1}$  then  $\eta_{i+1}$  is immediately a sibling of  $\eta_i$ . Otherwise the first transition is  $\circ_i \xrightarrow[\sigma_i]{\rightarrow} \bullet_i$ . The reached node is necessarily a sibling of  $\eta_i$ . To reach  $\eta_{i+1}$  only transitions labelled  $\rightarrow_+$  and  $\rightarrow$  can be followed. Hence, the node reached must also be a sibling of  $\eta_i$ .

Thus, by induction,  $\eta_n = \eta$  satisfies  $\sigma_1 o_1 \cdots o_{n-1} \sigma_n = \varphi$ .  $\square$

**Lemma C.2.** *For each CSS selector  $\varphi$  and tree  $T$ , we have*

$$T, \eta \models \varphi \Rightarrow (T, \eta) \in \mathcal{L}(\mathcal{A}_\varphi)$$

*Proof.* Assume  $T, \eta \models \varphi$ . Thus, since  $\varphi = \sigma_1 o_1 \cdots o_{n-1} \sigma_n$ , we have a sequence of nodes  $\eta_1, \dots, \eta_n$  such that for each  $i$  we have  $T, \eta_i \models \sigma_1 o_1 \cdots o_{i-1} \sigma_i$ . Note  $\eta_n = \eta$ . We build a run of  $\mathcal{A}_\varphi$  from  $\sigma_1$  to  $\circ_i$  by induction. When  $i = 1$  we have the run constructed by taking the loops on the initial state  $\circ_1$  labelled  $\downarrow$  and  $\rightarrow_+$  to navigate to  $\eta_1$ . Assume we have a run to  $\circ_i$ . We build a run to  $\circ_{i+1}$  we consider  $o_i$ .

- When  $o_i = \gg$  we know  $\eta_{i+1}$  is a descendant of  $\eta_i$ . We consider the construction of  $\mathcal{A}_\varphi$ . If  $\eta_{i+1}$  is the first child of  $\eta_i$ , we construct the run to  $\circ_{i+1}$  via the transition  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \circ_{i+1}$ , noting that we know  $\eta_i$  satisfies  $\sigma_i$ . Otherwise we take  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \bullet_i$  and arrive at either an ancestor or sibling of  $\eta_{i+1}$ . In the case of a neighbour, we can take the transition labelled  $\rightarrow$  to reach  $\eta_{i+1}$ . For an indirect sibling we can take the transition labelled  $\rightarrow_+$  followed by the transition labelled  $\rightarrow$ . For an ancestor, we take the transition labelled  $\downarrow$  and arrive at another sibling or ancestor of  $\eta_{i+1}$  that is closer. We continue in this way until we reach  $\eta_{i+1}$  as needed.
- When  $o_i = >$  we know  $\eta_{i+1}$  is a child of  $\eta_i$ . We consider the construction of  $\mathcal{A}_\varphi$ . If  $\eta_{i+1}$  is the first child of  $\eta_i$ , we construct the run to  $\circ_{i+1}$  via the transition  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \circ_{i+1}$ , noting that we know  $\eta_i$  satisfies  $\sigma_i$ . Otherwise we take  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \bullet_i$  and arrive at a preceding sibling of  $\eta_{i+1}$ . We can take the transition labelled  $\rightarrow_+$  to reach the preceding neighbour of  $\eta_{i+1}$  if required, and then the transition labelled  $\rightarrow$  to reach  $\eta_{i+1}$  as required.
- When  $o_i = +$  we know  $\eta_{i+1}$  is the neighbour of  $\eta_i$ . We consider the construction of  $\mathcal{A}_\varphi$  and take the only available transition  $\circ_i \xrightarrow[\sigma_i]{\rightarrow} \circ_{i+1}$ , noting that we know  $\eta_i$  satisfies  $\sigma_i$ . Thus, we reach  $\eta_{i+1}$  as required.
- When  $o_i = \sim$  we know  $\eta_{i+1}$  is a sibling of  $\eta_i$ . We consider the construction of  $\mathcal{A}_\varphi$ . If  $\eta_{i+1}$  is the neighbour of  $\eta_i$ , we construct the run to  $\circ_{i+1}$  via the transition  $\circ_i \xrightarrow[\sigma_i]{\downarrow} \circ_{i+1}$ , noting that we know  $\eta_i$  satisfies  $\sigma_i$ . Otherwise we take  $\circ_i \xrightarrow[\sigma_i]{\rightarrow} \bullet_i$  and arrive at a preceding sibling of  $\eta_{i+1}$ . We can

take the transition labelled  $\rightarrow_+$  to reach the preceding neighbour of  $\eta_{i+1}$  is required, and then the transition labelled  $\rightarrow$  to reach  $\eta_{i+1}$  as required.

Thus, by induction, we construct a run to  $\eta_n$  ending in state  $\circ_n$ . We transform this to an accepting run by taking the transition  $\circ_n \xrightarrow[\sigma_n]{\circ} q_f$ , using the fact that  $\eta_n$  satisfies  $\sigma_n$ .  $\square$

## C.2 Proof of Proposition 6.2

We show that

$$(T, \eta) \in \mathcal{L}(\mathcal{A}_1) \wedge (T, \eta) \in \mathcal{L}(\mathcal{A}_2) \Leftrightarrow (T, \eta) \in \mathcal{L}(\mathcal{A}_1 \cap \mathcal{A}_2).$$

We begin by observing that that all runs of a CSS automaton showing acceptance of a node  $\eta$  in  $T$  must follow a sequence of nodes  $\eta_1, \dots, \eta_n$  such that

- $\eta_1$  is the root of  $T$ , and
- when  $\eta_j = \eta' \iota$  then either  $\eta_{j+1} = \eta'(\iota + 1)$  or  $\eta_{j+1} = \eta_j 1$  for all  $j$ , and
- $\eta_n = \eta$

that defines the path taken by the automaton. Each node is “read” by some transition on each run. Note a transition labelled  $\rightarrow_+$  may read sequence nodes that is a factor of the path above. However, since these transitions are loops that do not check the nodes, without loss of generality we can assume each  $\rightarrow_+$  in fact reads only a single node. That is,  $\rightarrow_+$  behaves like  $\rightarrow$ . Recall,  $\rightarrow_+$  was only introduced to ensure the existence of “short” runs.

Because of the above, any two runs accepting  $\eta$  in  $T$  must follow the same sequence of nodes and be of the same length.

We have  $(T, \eta) \in \mathcal{L}(\mathcal{A}_1) \wedge (T, \eta) \in \mathcal{L}(\mathcal{A}_2)$  iff there are accepting runs

$$q_1^i \xrightarrow[\sigma_1^i]{d_1^i} \dots \xrightarrow[\sigma_n^i]{d_n^i} q_{n+1}^i$$

of  $\mathcal{A}_i$  over  $T$  reaching node  $\eta$  for both  $i \in \{1, 2\}$ . We argue these two runs exist iff we have a run

$$(q_1^1, q_1^2) \xrightarrow[\sigma_1]{d_1} \dots \xrightarrow[\sigma_n]{d_n} (q_{n+1}^1, q_{n+1}^2)$$

of  $\mathcal{A}_1 \cap \mathcal{A}_2$  where each  $d_j$  and  $\sigma_j$  depends on  $(d_j^1, d_j^2)$ .

- When  $(\downarrow, \downarrow)$  we have  $d_j = \downarrow$  and  $\sigma_j = \sigma_j^1 \cap \sigma_j^2$ .
- When  $(\rightarrow, \rightarrow)$  we have  $d_j = \rightarrow$  and  $\sigma_j = \sigma_j^1 \cap \sigma_j^2$ .
- When  $(\rightarrow, \rightarrow_+)$  we have  $d_j = \rightarrow$  and  $\sigma_j = \sigma_j^1$ .
- When  $(\rightarrow_+, \rightarrow)$  we have  $d_j = \rightarrow$  and  $\sigma_j = \sigma_j^2$ .
- When  $(\rightarrow_+, \rightarrow_+)$  we have  $d_j = \rightarrow_+$  and  $\sigma_j = *$ .
- When  $(\circ, \circ)$  we have  $d_j = \circ$  and  $\sigma_j = \sigma_j^1 \cap \sigma_j^2$ .
- The cases  $(\downarrow, \rightarrow_+)$ ,  $(\downarrow, \rightarrow)$ ,  $(\downarrow, \circ)$ ,  $(\rightarrow, \downarrow)$ ,  $(\rightarrow, \circ)$ ,  $(\rightarrow_+, \downarrow)$ ,  $(\rightarrow_+, \circ)$ ,  $(\circ, \downarrow)$ ,  $(\circ, \rightarrow)$ , and  $(\circ, \rightarrow_+)$  are not possible.

The existence of the transitions comes from the definition of  $\mathcal{A}_1 \cap \mathcal{A}_2$ . We have to argue that  $\eta_j$  satisfies both  $\sigma_j^i$  iff it also satisfies  $\sigma_j$ . By observing  $\sigma \cap * = * \cap \sigma = \sigma$  we always have  $\sigma_j = \sigma_j^1 \cap \sigma_j^2$ .

Let  $\sigma_j^i = \tau_i \Theta_i$  and  $\sigma_j = \tau \Theta$ . It is immediate that  $\eta_j$  satisfies  $\Theta = \Theta_1 \cup \Theta_2$  iff it satisfies both  $\Theta_i$ .

To complete the proof we need to show  $\eta_j$  satisfies  $\tau$  iff it satisfies both  $\tau_i$ . Note, we must have some  $s$  and  $e$  such that  $\tau, \tau_1, \tau_2 \in \{*, (s|*), (s|e), e\}$  else the type selectors cannot be satisfied (either  $\tau = \text{not } *$  or  $\tau_1$  and  $\tau_2$  assert conflicting namespaces or elements).

If some  $\tau_i = *$  the property follows by definition. Otherwise, if  $\tau = \tau_2$  then in all cases the conjunction of  $\tau_1$  and  $\tau_2$  is equivalent to  $\tau_2$  and we are done. The situation is similar when  $\tau = \tau_1$ . Otherwise  $\tau = (s|e)$  and  $\tau_1 = (s|*)$  and  $\tau_2 = e$  or vice versa, and it is easy to see  $\tau$  is equivalent to the intersection of  $\tau_1$  and  $\tau_2$ . Thus, we are done.

### C.3 Proofs for Non-Emptiness of CSS Automata

#### C.3.1 Bounding Namespaces and Elements

We show Proposition 6.4 (Bounded Types). We need to define the finite sets  $\downarrow(\text{ELE})$  and  $\downarrow(s)$ . To this end, we write

1.  $\text{ELE}_{\mathcal{A}}$  to denote the set of namespaced elements  $s:e$  such that there is some transition  $q \xrightarrow[\sigma]{d} q' \in \Delta$  with  $\sigma = (s|e) \Theta$  for some  $s, e$ , and  $\Theta$ ,
2.  $S_{\mathcal{A}}$  is the set of transitions  $q \xrightarrow[\sigma]{d} q' \in \Delta$  with  $\sigma \neq *$  and  $|\mathcal{A}|_{\sigma}$  denotes the cardinality of  $S_{\mathcal{A}}$ .

Let  $\{\tau_1, \dots, \tau_{|\mathcal{A}|_{\sigma}}\}$  be a set of fresh namespaced elements and

$$\downarrow(\text{ELE}_{\mathcal{A}}) = \text{ELE}_{\mathcal{A}} \uplus \{\tau_1, \dots, \tau_{|\mathcal{A}|_{\sigma}}\} \uplus \{\perp\}$$

where there is a bijection  $\theta : S_{\mathcal{A}} \rightarrow \{\tau_1, \dots, \tau_{|\mathcal{A}|_{\sigma}}\}$  such that for each  $t \in S_{\mathcal{A}}$  we have  $\theta(t) = \tau$  and

1.  $\tau = s:e$  if  $\sigma$  can only match elements  $s:e$ ,
2.  $\tau = s:e$  for some fresh element  $e$  if  $\sigma$  can only match elements of the form  $s:e'$  for all elements  $e'$ , and
3.  $\tau = s:e$  for some fresh namespace  $s$  if  $\sigma$  can only match elements of the form  $s':e$  for all namespaces  $s'$ , and
4.  $\tau = s:e$  for fresh  $s$  and fresh  $e$  if  $\sigma$  places no restrictions on the element type.

Thus, we can define bounded sets of namespaces and elements

$$\begin{aligned} \downarrow(\text{ELE}) &= \{e \mid \exists s. s:e \in \downarrow(\text{ELE}_{\mathcal{A}})\} \\ \downarrow(\text{NS}) &= \{s \mid \exists e. s:e \in \downarrow(\text{ELE}_{\mathcal{A}})\}. \end{aligned}$$

It remains to show  $\downarrow(\text{ELE}_{\mathcal{A}})$  is sufficient. That is, if some tree  $T$  is accepted by  $\mathcal{A}$ , we can define another tree  $T'$  that also is accepted by  $\mathcal{A}$  but only uses types in  $\downarrow(\text{ELE}_{\mathcal{A}})$ .

We take  $(T, \eta) \in \mathcal{L}(\mathcal{A})$  with  $T = (D, \lambda)$  and we define  $T' = (D, \lambda')$  satisfying the proposition. Let

$$q_0, \eta_0, q_1, \eta_1, \dots, q_{\ell}, \eta_{\ell}, q_{\ell+1}$$

be the accepting run of  $\mathcal{A}$ , by the sequence of transitions  $t_0, \dots, t_{\ell}$ . As noted above, we can assume each transition in  $\Delta$  appears only once in this sequence. Let  $\{\sigma_1, \dots, \sigma_{|\mathcal{A}|_{\sigma}}\}$  be the set of selectors appearing in  $\mathcal{A}$ . We perform the following modifications to  $\lambda$  to obtain  $\lambda'$ .



We obtain  $\lambda'$  from  $\lambda$  by changing the element labelling. We first consider all  $0 \leq i \leq \ell$  such that  $\eta_i$  is labelled by some element  $s:e \in \text{ELE}_{\mathcal{A}}$ . Let  $\text{Nodes}_{s:e}$  be the set of nodes labelled by  $s:e$  in  $\lambda$ . In  $\lambda'$  we label all nodes in  $\text{Nodes}_{s:e}$  by  $s:e$ . That is, we do not relabel nodes labelled by  $s:e$ . Let  $\text{Nodes}$  be the union of all such  $\text{Nodes}_{s:e}$ .

Next we consider all  $0 \leq i \leq \ell$  such that  $\eta_i \notin \text{Nodes}$  (i.e. was not labelled in the previous case) and  $t_i = q_i \xrightarrow{\sigma} q_{i+1}$  with  $\sigma \neq *$ . Let  $s:e \notin \text{ELE}_{\mathcal{A}}$  be the element labelling of  $\eta_i$  in  $\lambda$ . Moreover, take  $\tau$  such that  $\theta(t_i) = \tau$ . In  $\lambda'$  we label all nodes in  $\text{Nodes}_{s:e}$  (i.e. labelled by  $s:e$ ) in  $\lambda$  by  $\tau$ . That is, we globally replace  $s:e$  by  $\tau$ . Let  $\text{Nodes}'$  be  $\text{Nodes}$  union all such  $\text{Nodes}_e$ .

Finally, we label all nodes not in  $\text{Nodes}'$  with the null element  $\perp$ .

To see that

$$q_0, \eta_0, q_1, \eta_1, \dots, q_\ell, \eta_\ell, q_{\ell+1}$$

via  $t_0, \dots, t_\ell$  is an accepting run of  $(D, \lambda')$  we only need to show that for each  $t_i = q_i \xrightarrow{\sigma} q_{i+1}$  that  $\eta_i$  satisfies  $\sigma$ . This can be shown by induction over  $\sigma$ . Most atomic cases are straightforward (e.g. the truth of `:hover` is not affected by our transformations). The case of  $e$ ,  $(s|*)$ , or  $(s|e)$  appearing positively follows since in these cases the labelling remained the same or changed to some  $\tau$  consistent with the selector. When such selectors appear negatively, the result follows since we only changed elements and namespaces to fresh ones. The truth of attribute selectors remains unchanged since we did not change the attribute labelling. The cases of `:nth-child( $\alpha n + \beta$ )` and `:nth-last-child( $\alpha n + \beta$ )` follow since we did not change the number of nodes. For the selectors `:nth-of-type( $\alpha n + \beta$ )` and `:nth-last-of-type( $\alpha n + \beta$ )` there are two cases. If we did not change the element label  $s:e$  of  $\eta_i$ , then we also did not change the label of its siblings. Moreover, we did not add any  $s:e$  labels elsewhere in the tree. Hence the truth of the formulas remains the same. If we did change the label from  $s:e$  to  $\tau$  for some  $\tau$  then observe that we also relabelled all other nodes in the tree labelled by  $s:e$ . In particular, all siblings of  $\eta_i$ . Moreover, since  $\theta$  is a bijection and each transition appears only once in the run, we did not label any node not labelled  $s:e$  with  $\tau$ . Hence the truth of the formulas also remains the same. Similar arguments hold for `:only-child` and `:only-of-type`.

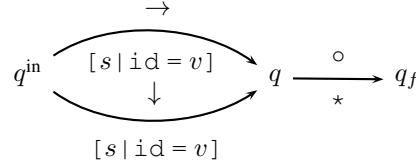
Thus,  $(D, \lambda')$  is accepted, and only uses elements in  $\downarrow(\text{ELE}_{\mathcal{A}})$  as required.

### C.3.2 Proof of Polynomial Bound on Attribute Value Lengths

We prove Proposition 6.5 (Bounded Attributes). That is we argue the existence of a polynomial bound for the solutions to any finite set  $C$  of constraints of the form  $[s|a \text{ op } v]$  or  $:\text{not}([s|a \text{ op } v])$ , for some fixed  $s$  and  $a$ . We say that  $C$  is a set of constraints over  $s$  and  $a$ .

In fact, the situation is a little more complicated because it may be the case that  $a$  is `id`. In this case we need to be able to enforce a global uniqueness constraint on the attribute values. Thus, for constraints on an ID attribute, we need a bound that is large enough to allow to all constraints on the same ID appearing throughout the automaton to be satisfied by unique values. Thus, for a given automaton, we might ask for a bound  $N$  such that if there exists unique ID values for each transition, then there exist values of length bounded by  $N$ .

However, the bound on the length must still work when we account for the fact that not all transitions in the automaton will be used during a run. Consider the following illustrative example.



In this case we have two transitions with ID constraints, and hence two sets of constraints  $C_1 = C_2 = \{[s \mid \text{id} = v]\}$ . Since these two sets of constraints cannot be satisfied simultaneously with unique values, even the bound  $N = 0$  will satisfy our naive formulation of the required property (since the property had the existence of a solution as an antecedent). However, it is easy to see that any run of the automaton does not use both sets of constraints, and that the bound  $N = |v|$  should suffice. Hence, we formulate the property of our bound to hold for all *sub-collections* of the collection of sets of constraints appearing in the automaton.

**Lemma C.3** (Bounded Attribute Values). *Given a collection of constraints  $C_1, \dots, C_n$  over some  $s$  and  $a$ , there exists a bound  $N$  polynomial in the size of  $C_1, \dots, C_n$  such that for any subsequence  $C_{i_1}, \dots, C_{i_m}$  if there is a sequence of words  $v_1, \dots, v_m$  such that all  $v_j$  are unique and  $v_j$  satisfies the constraints in  $C_{i_j}$ , then there is a sequence of words such that the length of each  $v_j$  is bounded by  $N$ , all  $v_j$  are unique, and  $v_j$  satisfies the constraints in  $C_{i_j}$ .*

The proof uses ideas from Muscholl and Walukiewicz’s NP fragment of LTL [45]. We first, for each set of constraints  $C$ , construct a deterministic finite word automaton  $\mathcal{A}$  that accepts only words satisfying all constraints in  $C$ . This automaton has a polynomial number of states and can easily be seen to have a short solution by a standard pumping argument. Given automata  $\mathcal{A}_1, \dots, \mathcal{A}_n$  with at most  $N_s$  states and  $N_c$  constraints in each set of constraints, we can again use pumping to show there is a sequence of distinct words  $v_1, \dots, v_n$  such that each  $v_i$  is accepted by  $\mathcal{A}_i$  and the length of  $v_i$  is at most  $n \cdot N_s \cdot N_c$ .

**The Automata** We define a type of word automata based on a model by Muscholl and Walukiewicz to show and NP upper bound for a variant of LTL. These automata read words and keep track of which constraints in  $C$  have been satisfied or violated. They accept once all positive constraints have been satisfied and no negative constraints have been observed.

In the following, let  $\text{Pref}(C)$  be the set of words  $v'$  such that  $v'$  is a prefix of some  $v$  with  $[s \mid a \text{ op } v] \in C$  or  $:\text{not}([s \mid a \text{ op } v]) \in C$ . Moreover, let  $\hat{\ } and  $\$$  be characters not in  $\Gamma$  that will mark the beginning and end of the word respectively. Additionally, let  $\varepsilon$  denote the empty word. Finally, we write  $v \preceq v'$  if  $v$  is a *factor* of  $v'$ , i.e.,  $v' = v_1 v v_2$  for some  $v_1$  and  $v_2$ .$

**Definition C.4** ( $\mathcal{A}_C$ ). *Given a set  $C$  of constraints over  $s$  and  $a$ , we define  $\mathcal{A}_C = (Q, \Delta, C)$  where*

- $Q$  is the set of all words  $a_1 v a_2$  such that
  - $v \in \text{Pref}(C)$ , and
  - $a_1, a_2 \in \Gamma \cup \{\varepsilon, \hat{\ }, \$\}$ .
- $\Delta \subseteq Q \times (\Gamma \cup \{\hat{\ }, \$\}) \times Q$  is the set of transitions  $v \xrightarrow{a} v'$  where  $v'$  is the longest suffix of  $va$  such that  $v' \in Q$ .

Observe that the size of the automaton  $\mathcal{A}_C$  is polynomial in the size of  $C$ .

A run of  $\mathcal{A}_C$  over a word with beginning and end marked  $a_1 \dots a_n \in \hat{\Gamma}^*\$$  is

$$(v_0, S_0, V_0) \xrightarrow{a_1} (v_1, S_1, V_1) \xrightarrow{a_2} \dots \xrightarrow{a_n} (v_n, S_n, V_n)$$

where  $v_0 = \varepsilon$  and for all  $1 \leq i \leq n$  we have  $v_{i-1} \xrightarrow{a_i} v_i$  and  $S_i, V_i \subseteq C$  track the satisfied and violated constraints respectively. That is  $S_0 = V_0 = \emptyset$ , and for all  $1 \leq i \leq n$  we have (noting  $\hat{v} \preceq v_i$  implies  $\hat{v}$  is a prefix of  $v_i$ , and similar for  $v\$$ )  $S_i =$

$$\begin{aligned} & S_{i-1} \cup \{[s|a = v] \in C \mid \hat{v}\$ = v_i\} \cup \\ & \{[s|a \sim = v] \in C \mid \exists a_1 \in \{\hat{\cdot}, \sqcup\}, a_2 \in \{\sqcup, \$\} . a_1 v a_2 \preceq v_i\} \cup \\ & \{[s|a \mid = v] \in C \mid \exists a_2 \in \{\$, -\} . \hat{v} a_2 \preceq v_i\} \cup \\ & \{[s|a \hat{=} v] \in C \mid \hat{v} \preceq v_i\} \cup \{[s|a \$ = v] \in C \mid v\$ \preceq v_i\} \cup \\ & \{[s|a * = v] \in C \mid v \preceq v_i\} \end{aligned}$$

and  $V_i =$

$$\begin{aligned} & V_{i-1} \cup \{:\text{not}([s|a = v]) \in C \mid \hat{v}\$ = v_i\} \cup \\ & \left\{ :\text{not}([s|a \sim = v]) \in C \mid \begin{array}{l} \exists a_1 \in \{\hat{\cdot}, \sqcup\}, \\ a_2 \in \{\sqcup, \$\} . a_1 v a_2 \preceq v_i \end{array} \right\} \cup \\ & \{:\text{not}([s|a \mid = v]) \in C \mid \exists a_2 \in \{\$, -\} . \hat{v} a_2 \preceq v_i\} \cup \\ & \{:\text{not}([s|a \hat{=} v]) \in C \mid \hat{v} \preceq v_i\} \cup \\ & \{:\text{not}([s|a \$ = v]) \in C \mid v\$ \preceq v_i\} \cup \\ & \{:\text{not}([s|a * = v]) \in C \mid v \preceq v_i\} . \end{aligned}$$

Such a run is *accepting* if  $S_n = \{[s|a \text{ op } v] \mid [s|a \text{ op } v] \in C\}$  and  $V_n = \emptyset$ . That is, all positive constraints have been satisfied and no negative constraints have been violated.

**Short Solutions** We show the existence of short solutions via the following lemma. The proof of this lemma is a simple pumping argument which appears below. Intuitively, if a satisfying word is shorter than  $N_s \cdot N_c$  we do not change it. If it is longer than  $N_s \cdot N_c$  any accepting run of the automaton on this word must contain a repeated  $(v, S, V)$ . We can thus pump down this word to ensure that it is shorter than  $N_s \cdot N_c$ . Then, to ensure it is unique, we pump it up to some unique length of at most  $n \cdot N_s \cdot N_c$ .

**Lemma C.5** (Short Attribute Values). *Given a sequence of sets of constraint automata  $\mathcal{A}_{C_1}, \dots, \mathcal{A}_{C_n}$  each with at most  $N_s$  states and at most  $N_c$  constraints in each  $C_i$ , if there is a sequence of pairwise unique words  $v_1, \dots, v_n$  such that for all  $1 \leq i \leq n$  there is an accepting run of  $\mathcal{A}_{C_i}$  over  $v_i$ , then there exists such a sequence where the length of each  $v_i$  is at most  $n \cdot N_s \cdot N_c$ .*

To obtain Lemma C.3 (Bounded Attribute Values) we observe that for any subsequence  $C_{i_1}, \dots, C_{i_m}$  we have  $m \cdot N'_s \cdot N'_c \leq n \cdot N_s \cdot N_c$  since  $m \leq n$  and the max number of states  $N'_s$  and constraints  $N'_c$  in the subsequence have  $N'_s \leq N_s$  and  $N'_c \leq N_c$ .

We give the proof of Lemma C.5. That is, given a sequence of sets of constraint automata  $\mathcal{A}_{C_1}, \dots, \mathcal{A}_{C_n}$  each with at most  $N_s$  states and at most  $N_c$  constraints in each  $C_i$ , if there is a sequence of pairwise unique words  $v_1, \dots, v_n$  such that for all  $1 \leq i \leq n$  there is an accepting run of  $\mathcal{A}_{C_i}$  over  $v_i$ , then there exists such a sequence where the length of each  $v_i$  is at most  $n \cdot N_s \cdot N_c$ .

To prove the lemma, take a sequence  $v_1, \dots, v_n$  such that each  $v_i$  is unique and accepted by  $\mathcal{A}_{C_i}$ . We proceed by induction, constructing  $v'_1, \dots, v'_i$  such that each  $v'_j$  is unique, accepted by  $\mathcal{A}_{C_j}$ , and of length  $\ell$  such that either

- $\ell \leq N_s \cdot N_c$  and  $v'_j = v_j$ , or
- $i \cdot N_s \cdot N_c \leq \ell \leq (i+1) \cdot N_s \cdot N_c$ .

When  $i = 0$  the result is vacuous. For the induction there are two cases.

When the length  $\ell$  of  $v_i$  is such that  $\ell \leq N_s \cdot N_c$  we set  $v'_i = v_i$ . We know  $v'_i$  is unique amongst  $v'_1, \dots, v'_i$  since for all  $j < i$  either  $v'_j$  is longer than  $v'_i$  or is equal to  $v_j$  and thus distinct from  $v_i$ .

When  $\ell > N_s \cdot N_c$  we use a pumping argument to pick some  $v'_i$  of length  $\ell'$  such that  $i \cdot N_s \cdot N_c \leq \ell' \leq (i+1) \cdot N_s \cdot N_c$ . This ensures that  $v'_i$  is unique since it is the only word whose length lies within the bound. We take the accepting run

$$(u_0, S_0, V_0) \xrightarrow{a_1} (u_1, S_1, V_1) \xrightarrow{a_2} \dots \xrightarrow{a_n} (u_\ell, S_n, V_\ell)$$

of  $v_i$  and observe that the values of  $S_j$  and  $V_j$  are increasing by definition. That is  $S_j \subseteq S_{j+1}$  and  $V_j \subseteq V_{j+1}$ . By a standard down pumping argument, we can construct a short accepting run containing only distinct configurations of length bound by  $N_s \cdot N_c$ . We construct this run by removing all cycles from the original run. This maintains the acceptance condition. Next we obtain an accepted word of length  $i \cdot N_s \cdot N_c \leq \ell' \leq (i+1) \cdot N_s \cdot N_c$ . Since  $\ell > N_s \cdot N_c$  we know there exists at least one configuration  $(u, S, V)$  in the short run that appeared twice in the original run. Thus there is a run of the automaton from  $(u, S, V)$  back to  $(u, S, V)$  which can be bounded by  $N_s \cdot N_c$  by the same downward pumping argument as before. Thus, we insert this run into the short run the required number of times to obtain an accepted word  $v'_i$  of the required length.

Thus, by induction, we are able to obtain the required short words  $v'_1, \dots, v'_n$  as needed.

### C.3.3 Missing definitions for $\text{AttsPres}(\theta, \vec{x})$

$$\begin{aligned} \text{AttsPres}([s \mid a \sim = v], \vec{x}) &= \left( \bigwedge_{1 \leq j \leq m} \left( \begin{array}{c} x_{i,j}^{s:a} = a_j \wedge \\ x_{i,m+1}^{s:a} = 0 \\ \vee \\ x_{i,m+1}^{s:a} = \perp \end{array} \right) \right) \\ \text{AttsPres}([s \mid a \mid = v], \vec{x}) &= \bigwedge_{1 \leq j \leq m} x_{i,j}^{s:a} = a_j \wedge \left( \begin{array}{c} x_{i,j-1}^{s:a} = \perp \wedge \\ \bigwedge_{1 \leq j' \leq m} x_{i,j+j'}^{s:a} = a_j \wedge \\ (x_{i,j+m+1}^{s:a} = 0 \vee x_{i,j+m+1}^{s:a} = \perp) \end{array} \right) \\ \text{AttsPres}([s \mid a \S = v], \vec{x}) &= \bigvee_{0 \leq j \leq N-m-1} \left( \begin{array}{c} \bigwedge_{1 \leq j' \leq m} x_{i,j+j'}^{s:a} = a_j \wedge \\ x_{i,j+m+1}^{s:a} = 0 \end{array} \right) \end{aligned}$$

### C.3.4 Negating Positional Formulas

We need to negate selectors like  $\text{:nth-child}(\alpha n + \beta)$ , For completeness, we give the definition of the negation below.

We decompose  $\beta$  according to the period  $\alpha$ . I.e.  $\beta = \alpha\beta_1 + \beta_2$ , where  $\beta_1$  and  $\beta_2$  are the unique integers such that  $|\beta_2| < |\alpha|$  and  $\beta_1\alpha < 0$  implies  $\beta_2 \leq 0$  and  $\beta_1\alpha > 0$  implies  $\beta_2 \geq 0$ .

**Definition C.6** ( $\text{NoMatch}(\bar{x}, \alpha, \beta)$ ). *Given constants  $\alpha, \beta, \beta_1$ , and  $\beta_2$  as above, we define the formula*

$\text{NoMatch}(\bar{x}, \alpha, \beta)$  to be

$$(0 \geq \alpha \wedge \bar{x} < \beta) \vee (0 \leq \alpha \wedge \bar{x} > \beta) \vee \left( \exists \bar{n}. \exists \bar{\beta}'_2. \left( \begin{array}{l} |\bar{\beta}'_2| < |\alpha| \wedge \\ \left( \beta_1 \alpha < 0 \Rightarrow \bar{\beta}'_2 \leq 0 \right) \wedge \\ \left( \beta_1 \alpha > 0 \Rightarrow \bar{\beta}'_2 \geq 0 \right) \wedge \\ \bar{\beta}'_2 \neq \beta_2 \wedge \\ \bar{x} = \alpha \bar{n} + \alpha \beta_1 + \bar{\beta}'_2 \end{array} \right) \right)$$

In the following, whenever we negate a formula of the form  $\neg(\exists \bar{n}. \bar{x} = \alpha \bar{n} + \beta)$  we will use the formula  $\text{NoMatch}(\bar{x}, \alpha, \beta)$ . One can verify that the resulting formula is existential Presburger. We show that our negation of periodic constraints is correct.

**Proposition C.7** (Correctness of  $\text{NoMatch}(\bar{x}, \alpha, \beta)$ ). *Given constants  $\alpha$  and  $\beta$ , we have*

$$\neg(\exists \bar{n}. \bar{x} = \alpha \bar{n} + \beta) \Leftrightarrow \text{NoMatch}(\bar{x}, \alpha, \beta).$$

*Proof.* We first consider  $\alpha = 0$ . Since there is no  $\beta'_2$  with  $|\beta'_2| < 0$  we have to prove

$$\neg(\bar{x} = \beta) \Leftrightarrow (\bar{x} < \beta) \vee (\bar{x} > \beta)$$

which is immediate.

In all other cases, the conditions on  $\beta_2$  and  $\beta'_2$  ensure that we always have  $0 < |\beta_2 - \beta'_2| < |\alpha|$ .

If  $\exists \bar{n}. \bar{x} = \alpha \bar{n} + \alpha \beta_1 + \beta_2$  then if  $\alpha > 0$  it is easy to verify that we don't have  $\bar{x} < \beta$  (since  $\bar{n} = 0$  gives  $\bar{x} = \beta$  as the smallest value of  $\bar{x}$ ) and similarly when  $\alpha < 0$  we don't have  $\bar{x} < \beta$ . To disprove the final disjunct of  $\text{NoMatch}(\bar{x}, \alpha, \beta)$  we observe there can be no  $\bar{n}'$  s.t.  $\bar{x} = \alpha \bar{n}' + \alpha \beta_1 + \beta'_2$  since  $\bar{x} = \alpha \bar{n} + \alpha \beta_1 + \beta_2$  and  $0 < |\beta_2 - \beta'_2| < |\alpha|$ .

In the other direction, we have three cases depending on the satisfied disjunct of  $\text{NoMatch}(\bar{x}, \alpha, \beta)$ . Consider  $\alpha > 0$  and  $\bar{x} < \beta$ . In this case there is no  $\bar{n}$  such that  $\bar{x} = \alpha \bar{n} + \alpha \beta_1 + \beta_2 = \alpha \bar{n} + \beta$  since  $\bar{n} = 0$  gives the smallest value of  $\bar{x}$ , which is  $\beta$ . The case is similar for the second disjunct with  $\alpha < 0$ .

The final disjunct gives some  $\bar{n}'$  such that  $\bar{x} = \alpha \bar{n}' + \alpha \beta_1 + \beta'_2$  with  $0 < |\beta_2 - \beta'_2| < |\alpha|$ . Hence, there can be no  $\bar{n}$  with  $\bar{x} = \alpha \bar{n} + \alpha \beta_1 + \beta_2$ .  $\square$

### C.3.5 Correctness of Presburger Encoding

We prove soundness and completeness of the Presburger encoding of CSS automata non-emptiness in the two lemmas below.

**Lemma C.8.** *For a CSS automaton  $\mathcal{A}$ , we have*

$$\mathcal{L}(\mathcal{A}) \neq \emptyset \Rightarrow \theta_{\mathcal{A}} \text{ is satisfiable.}$$

*Proof.* We take a run of  $\mathcal{A}$  and construct a satisfying assignment to the variables in  $\theta_{\mathcal{A}}$ . That is take a document tree  $T = (D, \lambda)$ , node  $\eta \in D$ , and sequence

$$q_0, \eta_0, q_1, \eta_1, \dots, q_\ell, \eta_\ell, q_{\ell+1} \in (Q \times D)^* \times \{q_f\}$$

that is an accepting run. We know from Proposition 6.4 (Bounded Types) that  $T$  can only use namespaces from  $\downarrow(\text{NS})$  and elements from  $\downarrow(\text{ELE})$ . Let  $t_0, \dots, t_\ell$  be the sequence of transitions used in the accepting run. We assume (w.l.o.g.) that no transition is used twice. We construct a satisfying assignment to the variables as follows.

- $\bar{q}_i = q_i$  for all  $i \leq \ell + 1$  and  $\bar{q}_i = q_f$  for all  $i > \ell + 1$ .
- $\bar{s}_i = \lambda_S(\eta_i)$  for all  $i \leq \ell + 1$  ( $\bar{s}_i$  can take any value for other values of  $i$ ).
- $\bar{e}_i = \lambda_E(\eta_i)$  for all  $i \leq \ell + 1$  ( $\bar{e}_i$  can take any value for other values of  $i$ ).
- $\bar{p}_i = (p \in \lambda_P(\eta_i))$ , for each pseudo-class  $p \in P \setminus \{\text{root}\}$  and  $i \leq \ell + 1$  (these variables can take any value for  $i > \ell + 1$ ).
- $\bar{n}_i = \iota$ , when  $\eta_i = \eta' \iota$  for some  $\eta'$  and  $\iota$  and  $1 \leq i \leq \ell + 1$ , otherwise  $\bar{n}_i$  can take any value.
- $\bar{n}_i^{s:e} = j$ , where  $j$  is the number of nodes of type  $s:e$  preceding  $\eta_i$  in the sibling order. That is  $\eta_i = \eta' \iota$  for some  $\eta'$  and  $\iota$  and

$$j = \left| \left\{ \eta' \iota' \mid \begin{array}{l} \iota' < \iota \wedge \eta' \iota' \in D \wedge \\ \lambda_S(\eta' \iota') = \lambda_S(\eta) \wedge \\ \lambda_E(\eta' \iota') = \lambda_E(\eta) \end{array} \right\} \right|.$$

When  $i = 0$  or  $i > \ell + 1$  we can assign any value to  $\bar{n}_i^{s:e}$ .

- $\bar{N}_i = N - \iota$  where  $i \leq \ell + 1$  and  $\eta = \eta' \iota$  for some  $\eta'$  and  $\iota$  and  $N$  is the smallest number such that  $\eta' N \notin D$ . For  $i = 0$  or  $i > \ell + 1$  the variable  $\bar{N}_i$  can take any value.
- $\bar{N}_i^{s:e} = j$ , where  $j$  is the number of nodes of type  $s:e$  succeeding  $\eta_i$  in the sibling order. That is  $\eta_i = \eta' \iota$  for some  $\eta'$  and  $\iota$  and

$$j = \left| \left\{ \eta' \iota' \mid \begin{array}{l} \iota' > \iota \wedge \eta' \iota' \in D \wedge \\ \lambda_S(\eta' \iota') = \lambda_S(\eta) \wedge \\ \lambda_E(\eta' \iota') = \lambda_E(\eta) \end{array} \right\} \right|.$$

When  $i = 0$  or  $i > \ell + 1$  we can assign any value to  $\bar{N}_i^{s:e}$ .

- Assignments to  $x_{i,j}^{s:a}$  are discussed below.

It remains to prove that the given assignment satisfies the formula.

Recall

$$\theta_{\mathcal{A}} = \left( \begin{array}{l} \bar{q}_0 = q^{\text{in}} \wedge \bar{q}_n = q_f \wedge \\ \bigwedge_{0 \leq i < n} (\text{Tran}(i) \vee \bar{q}_i = q_f) \wedge \\ \text{Consistent} \end{array} \right).$$

The first two conjuncts follow immediately from our assignment to  $\bar{q}_i$  and that the chosen run was accepting. Next we look at the third conjunct and simultaneously prove  $\text{Consistent}_n$ . When  $i \geq \ell + 1$  we assigned  $q_f$  to  $\bar{q}_i$  and can choose any assignment that satisfies  $\text{Consistent}_n$ . Otherwise we show we satisfy  $\text{Tran}(i)$  by showing we satisfy  $\text{Tran}(i, t_i)$ . We also show  $\text{Consistent}_n$  is satisfied by induction, noting it is immediate for  $i = 0$  and that for  $i = 1$  we must have either the first or last case which do not depend on the induction hypothesis. Consider the form of  $t_i$ .

1. When  $t_i = q_i \xrightarrow{\downarrow \sigma} q_{i+1}$  we immediately confirm the values of  $\bar{q}_i, \bar{q}_{i+1}, \bar{n}_{i+1}, \bar{n}_{i+1}^{s:e}$  satisfy the constraint. Similarly for  $\neg : \text{empty}_i$  since we know  $: \text{empty} \notin \lambda_P(\eta_i)$ . We defer the argument for  $\text{Pres}(\sigma, i)$  until after the case split. That  $\text{Consistent}_n$  is satisfied can also be seen directly.

2. When  $t_i = q_i \xrightarrow{\sigma} q_{i+1}$  we know  $\eta_i = \eta'\iota$  and  $\eta_{i+1} = \eta'(\iota + 1)$  for some  $\eta'$  and  $\iota$ . We can easily check the values of  $\bar{q}_i, \bar{q}_{i+1}, \bar{n}_{i+1}, \bar{N}_i, \bar{n}_{i+1}^{s:e},$  and  $\bar{N}_{i+1}^{s:e}$  satisfy the constraint. We defer the argument for  $\text{Pres}(\sigma, i)$  until after the case split. To show  $\text{Consistent}_n$  we observe  $\bar{n}_{i+1}$  is increased by 1 and only one  $\bar{n}_{i+1}^{s:e}$  is increased by 1, the others being increased by 0. Similarly for  $\bar{N}_i$  and  $\bar{n}_{i+1}^{s:e}$ . Hence the result follows from induction.
3. When  $t_i = q_i \xrightarrow[\ast]{\sigma} q_{i+1}$  we know  $\eta_i = \eta'\iota$  and  $\eta_{i+1} = \eta'(\iota')$  for some  $\eta', \iota,$  and  $\iota < \iota'$ . We can easily check the values of  $\bar{q}_i, \bar{q}_{i+1}, \bar{n}_{i+1}, \bar{N}_i, \bar{n}_{i+1}^{s:e},$  and  $\bar{N}_{i+1}^{s:e}$  satisfy the constraint. We defer the argument for  $\text{Pres}(\sigma, i)$  until after the case split. To satisfy the constraints over the position variables, we observe that values for  $\bar{\delta}$  and  $\bar{\delta}_{s:e}$  can be chosen easily for the specified assignment. Combined with induction this shows  $\text{Consistent}_n$  as required.
4. When  $t_i = q_i \xrightarrow[\sigma]{\circ} q_{i+1}$

We can easily check the values of  $\bar{q}_i$  and  $\bar{q}_{i+1}$ . We defer the argument for  $\text{Pres}(\sigma, i)$  until after the case split. By induction we immediately obtain  $\text{Consistent}_n$ .

We show  $\text{Pres}(\sigma, i)$  is satisfied for each  $\eta_i$  and  $\sigma$  labelling  $t_i$ . Take a node  $\eta$  and  $\sigma = \tau\Theta$  from this sequence. Note  $\eta$  satisfies  $\sigma$  since the run is accepting. Recall

$$\text{Pres}(\tau\Theta, i) = \left( \begin{array}{c} \text{Pres}(\tau, i) \wedge \\ \left( \bigwedge_{\theta \in \text{NoAtts}(\Theta)} \text{Pres}(\theta, i) \right) \wedge \\ \text{AttsPres}(\tau\Theta, i) \end{array} \right).$$

From the type information of  $\eta$  we immediately satisfy  $\text{Pres}(\tau, i)$ .

For a positive  $\theta \in \Theta$  there are several cases. If  $\theta = \text{:root}$  then we know we are in  $\eta_0$  and the encoding is  $\top$ . If  $\theta$  is some other pseudo class  $p$  then the encoding of  $\theta$  is  $\bar{p}_i$  and we assigned true to this variable. For  $\text{:nth-child}(\alpha n + \beta)$  and  $\text{:nth-last-child}(\alpha n + \beta)$  satisfaction of the encoding follows immediately from  $\eta$  satisfying  $\theta$  and our assignment to  $\bar{n}_i$  and  $\bar{N}_i$ . We satisfy the encodings of  $\text{:nth-of-type}(\alpha n + \beta), \text{:nth-last-of-type}(\alpha n + \beta), \text{:only-child},$  and  $\text{:only-of-type}$  similarly. The latter follow since an only child is position 1 from the start and end, and an only of type node has 0 strict predecessors or successors of the same type.

For a negative  $\theta \in \Theta$  there are several cases. If  $\theta = \text{:not}(\text{:root})$  then we know we are not in  $\eta_0$  and the encoding is  $\top$ . If  $\theta$  is the negation of some other pseudo class  $p$  then the encoding of  $\theta$  is  $\neg\bar{p}_i$  and we assigned false to this variable. For the selectors  $\text{:not}(\text{:nth-child}(\alpha n + \beta))$  and the opposite selector  $\text{:not}(\text{:nth-last-child}(\alpha n + \beta))$  satisfaction of the encoding follows immediately from  $\eta$  satisfying  $\theta$ , our assignment to  $\bar{n}_i$  and  $\bar{N}_i$  as well as Proposition C.7 (Correctness of  $\text{NoMatch}(\bar{x}, \alpha, \beta)$ ). We satisfy encodings of  $\text{:not}(\text{:nth-of-type}(\alpha n + \beta))$  and of  $\text{:not}(\text{:nth-last-of-type}(\alpha n + \beta))$  in a likewise fashion. For the remaining cases of  $\text{:not}(\text{:only-child}),$  and  $\text{:not}(\text{:only-of-type})$  the property follows since, for the former the node must either not be position 1 from the start or end, and for the latter a not only of type node has more than 0 strict predecessors or successors of the same type.

Next, to satisfy  $\text{AttsPres}(\tau\Theta, i)$  we have to satisfy a number of conjuncts. First, if we have a word  $a_1 \dots a_n$  we assign it to the variables  $x_{i,j}^{s:a}$  (where  $\eta$  is the  $i$ th in the run and  $j$  ranges over all word positions within the computed bound) by assigning  $x_{i,j}^{s:a} = a_j$  when  $j \leq n$  and  $x_{i,j}^{s:a} = 0$  otherwise.

In all cases below, it is straightforward to observe that if a word (within the computed length bound) satisfies  $[s|a \text{ op } v]$  or  $\text{:not}([s|a \text{ op } v])$  then the encoding  $\text{AttsPres}_{s:a}([s|a \text{ op } v], i)$  or

$\neg \text{AttsPres}_{s:a}([s \mid a \text{ op } v], i)$  is satisfied by our variable assignment. Similarly  $\text{Nulls}(\vec{x})$  is straightforwardly satisfied. Hence, if a word satisfies  $C$  then our assignment to the variables means  $\text{AttsPres}_{s:a}(C, i)$  is also satisfied.

There are a number of cases of conjuncts for attribute selectors. The simplest is for sets  $\Theta_a^s$  where we see immediately that all constraints are satisfied for  $\lambda_{\mathbb{A}}(\eta)(s, a)$  and hence we assign this value to the appropriate variables and the conjunct is satisfied also. For each  $[a]$  and  $[a \text{ op } v] \in \Theta$  we have in the document some namespace  $s$  such that  $\lambda_{\mathbb{A}}(\eta)(s, a)$  satisfies the attribute selector and all negative selectors applying to all namespaces. Let  $s'$  be the fresh name space assigned to the selector during the encoding and  $C$  be the full set of constraints belonging to the conjunct (i.e. including negative ones). We assign to the variable  $x_{i,j}^{s:a}$  the  $j$ th character of  $\lambda_{\mathbb{A}}(\eta)(s, a)$  (where  $\eta$  is the  $i$ th in the run) and satisfy the conjunct as above. Note here that a single value of  $s:a$  is assigned to several  $s':a$ . This is benign with respect to the global uniqueness required by ID attributes because each copy has a different namespace.

Finally, we have to satisfy the consistency constraints. We showed  $\text{Consistent}_n$  above. The remaining consistency constraints are easily seen to be satisfied:  $\text{Consistent}_i$  because each ID is unique causing at least one pair of characters to differ in every value;  $\text{Consistent}_p$  since it encodes basic consistency constraints on the appearance of pseudo elements in the tree.

Thus, we have satisfied the encoded formula, completing the first direction of the proof.  $\square$

**Lemma C.9.** *For a CSS automaton  $\mathcal{A}$ , we have*

$$\theta_{\mathcal{A}} \text{ is satisfiable.} \Rightarrow \mathcal{L}(\mathcal{A}) \neq \emptyset$$

*Proof.* Take a satisfying assignment  $\rho$  to the free variables of  $\theta_{\mathcal{A}}$ . We construct a tree and node  $(T, \eta)$  as well as a run of  $\mathcal{A}$  accepting  $(T, \eta)$ .

We begin by taking the sequence of states  $q_0, \dots, q_{\ell+1}$  which is the prefix of the assignment to  $\bar{q}_0, \dots, \bar{q}_n$  where  $q_{\ell}$  is the first occurrence of  $q_f$ . We will construct a series of transitions  $t_0, \dots, t_{\ell}$  with  $t_i = q_i \xrightarrow[\sigma_i]{d_i} q_{i+1}$  for all  $0 \leq i \leq \ell$ . We will define each  $d_i$  and  $\sigma_i$ , as well as construct  $T$  and  $\eta$  by induction. We construct the tree inductively, then show  $\sigma_i$  is satisfied for each  $i$ .

At first let  $T_0$  contain only a root node. Thus  $\eta_0$  is necessarily this root. Throughout the proof we label each  $\eta_i$  as follows.

- $\lambda_{\mathbb{S}}(\eta_i) = \rho(\bar{s}_i)$  (i.e. we assign the value given to  $\bar{s}_i$  in the satisfying assignment).
- $\lambda_{\mathbb{E}}(\eta_i) = \rho(\bar{e}_i)$ .
- $\lambda_{\mathbb{P}}(\eta_i) = \left( \begin{array}{l} \{p \mid p \in P \setminus \{\text{:root}\} \wedge \rho(\bar{p}_i) = \top\} \cup \\ \{\text{:root} \mid i = 0\} \end{array} \right)$ .
- $\lambda_{\mathbb{A}}(\eta_i)(s, a) = \rho(x_{i,1}^{s:a} \dots x_{i,N}^{s:a})$  where  $\rho(x_{i,1}^{s:a} \dots x_{i,N}^{s:a})$  is the word obtained by stripping all of the null characters from  $\rho(x_{i,1}^{s:a}) \dots \rho(x_{i,N}^{s:a})$ .

We pick  $t_i$  as the transition corresponding to a satisfied disjunct of  $\text{Tran}(i)$  (of which there is at least one since  $q_i \neq q_f$  when  $i \leq \ell$ ). Thus, take  $t_i = q_i \xrightarrow[\sigma_i]{d_i} q_{i+1}$ . We proceed by a case split on  $d_i$ . Note only cases  $d_i = \downarrow$  and  $d_i = \circ$  may apply when  $i = 0$ .

- When  $d_i = \downarrow$  we build  $T_{i+1}$  as follows. First we add the leaf node  $\eta_{i+1} = \eta_i 1$ . Then, if  $i > 0$ , we add siblings appearing after  $\eta_i$  with types required by the last of type information. That is, we add  $\rho(\bar{N}_i) - 1 = \sum_{s:e \in \text{ELE}} \rho(\bar{N}_i^{s:e})$  siblings appearing after  $\eta_i$ . In particular, for each  $s$  and  $e$  we add  $\rho(\bar{N}_i^{s:e})$  new nodes. Letting  $\eta_i = \eta \nu$  each of these new nodes  $\eta'$  will have the form  $\eta \nu'$  with  $\nu' > \nu$ . We set  $\lambda_{\mathbb{S}}(\eta') = s$ ,  $\lambda_{\mathbb{E}}(\eta') = e$ ,  $\lambda_{\mathbb{P}}(\eta') = \emptyset$ ,  $\lambda_{\mathbb{A}}(\eta') = \emptyset$ .



- When  $d_i = \Rightarrow$  we build  $T_{i+1}$  by adding a single node to  $T_i$ . When  $\eta_i = \eta_\iota$  we add  $\eta_{i+1} = \eta(\iota + 1)$  with the labelling as above.
- When  $d_i = \Rightarrow_+$  we build  $T_{i+1}$  as follows. We add  $\rho(\bar{\delta}) = (\rho(\bar{n}_{i+1}) - \rho(\bar{n}_i))$  new nodes of the form  $\eta\iota'$  where  $\eta_i = \eta_\iota$  and  $\rho(\bar{n}_i) = \iota < \iota' \leq \rho(\bar{n}_i)$ . Let  $\eta_{i+1}$  be  $\eta\rho(\bar{n}_i)$  labelled as above. For the remaining new nodes, for each  $s$  and  $e$ , we label  $\rho(\bar{\delta}_{s:e})$  of the new nodes  $\eta'$  with  $\lambda_S(\eta') = s$ ,  $\lambda_E(\eta') = e$ ,  $\lambda_P(\eta') = \emptyset$ ,  $\lambda_A(\eta') = \emptyset$ . Note  $\text{Consistent}_n$  ensures we have enough new nodes to partition like this.
- When  $d_i = \circ$  and  $i = 0$  we have completed building the tree. If  $i > 0$ , we add siblings appearing after  $\eta_i$  with types required by the last of type information exactly as in the case of  $d = \Downarrow$  above.

The tree and node we require are the tree and node obtained after reaching some  $d_i = \circ$ , for which we necessarily have  $i = \ell$  since  $\circ$  must be and can only be used to reach  $q_f$ . In constructing this tree we have almost demonstrated an accepting run of  $\mathcal{A}$ . To complete the proof we need to argue that all  $\sigma_i$  are satisfied by  $\eta_i$  and that the obtained is valid. Let  $\tau\Theta = \sigma_i$ .

To check  $\tau$  we observe that  $\text{Pres}(\tau, i)$  constrains  $\bar{s}_i$  and  $\bar{e}_i$  to values, which when assigned to  $\eta_i$  as above mean  $\eta_i$  directly satisfies  $\tau$ .

Now, take some  $\theta \in \Theta$ . In each case we argue that  $\text{Pres}(\tau, i)$  ensures the needed properties. Note this is straightforward for the attribute selectors due to the directness of the Presburger encoding. Consider the remaining selectors.

First assume  $\theta$  is positive. If it is  $:\text{root}$  then we must have  $i = 0$  and  $\eta_i$  is the root node as required. For other pseudo classes  $p$  we asserted  $\bar{p}_i$  hence we have  $p \in \lambda_P(\eta_i)$ . The encoding of the remaining positive constraints can only be satisfied when  $i > 0$ . That is,  $\eta_i$  is not the root node.

For  $:\text{nth-child}(\alpha n + \beta)$  observe we constructed  $T$  such that  $\eta_i = \eta\rho(\bar{n}_i)$  for some  $\eta$ . From the defined encoding of  $\text{Pres}(\text{nth-child}(\alpha n + \beta), i)$  we directly obtain that  $\eta_i$  satisfies the selector  $:\text{nth-child}(\alpha n + \beta)$ . Similarly for  $:\text{nth-last-child}(\alpha n + \beta)$  as we always pad the end of the sibling order to ensure the correct number of succeeding siblings.

For  $:\text{nth-of-type}(\alpha n + \beta)$  and  $:\text{nth-last-of-type}(\alpha n + \beta)$  selectors, by similar arguments to the previous selectors, we have ensured that there are enough preceding or succeeding nodes (along with the directness of their Presburger encoding) to ensure these selectors are satisfied by  $\eta_i$  in  $T$ .

For  $:\text{only-child}$  we know there are no other children since  $\rho(\bar{n}_i) = \rho(\bar{N}_i) = 1$ . Finally for the selector  $:\text{only-of-type}$  we know there are no other children of the same type since  $\rho(\bar{n}_i^{s:e}) = \rho(\bar{N}_i^{s:e}) = 0$  where  $\eta_i$  has type  $s:e$ .

When  $\theta$  is negative there are several cases. If it is  $:\text{not}(\text{root})$  then we must have  $i > 0$  and  $\eta_i$  is not the root node. For other pseudo classes  $p$  we asserted  $\neg\bar{p}_i$  hence we have  $p \notin \lambda_P(\eta_i)$ . The encoding of the remaining positive constraints are always satisfied on the root node. That is,  $i = 0$ . When  $\eta_i$  is not the root node we have  $i > 0$ .

For  $:\text{not}(\text{nth-child}(\alpha n + \beta))$  observe we constructed  $T$  such that  $\eta_i = \eta\rho(\bar{n}_i)$  for some  $\eta$ . From the definition of  $\text{Pres}(\text{not}(\text{nth-child}(\alpha n + \beta)), i)$  we obtain that  $\eta_i$  does not satisfy the required selector  $:\text{nth-child}(\alpha n + \beta)$  via Proposition C.7 (Correctness of  $\text{NoMatch}(\bar{x}, \alpha, \beta)$ ). Similarly for the last child selector

$$:\text{not}(\text{nth-last-child}(\alpha n + \beta)).$$

For  $:\text{not}(\text{nth-of-type}(\alpha n + \beta))$  and  $:\text{not}(\text{nth-last-of-type}(\alpha n + \beta))$ , by similar arguments to the previous selectors, we have ensured that there are enough preceding or succeeding nodes (along with their Presburger encodings and Proposition C.7 (Correctness of  $\text{NoMatch}(\bar{x}, \alpha, \beta)$ )) to ensure these selectors are satisfied by  $\eta_i$  in  $T$ .

For `:not (:only-child)` we know there are some other children since  $\rho(\bar{n}_i) > 1$  or  $\rho(\bar{N}_i) > 1$ . Finally for `:not (:only-of-type)` we know there are other children of the same type since  $\rho(\bar{n}_i^{s:e}) > 0$  or  $\rho(\bar{N}_i^{s:e}) > 0$  where  $\eta_i$  has type  $s:e$ .

Thus we have an accepting run of  $\mathcal{A}$  over some  $(T, \eta)$ . However, we finally have to argue that  $T$  is a valid document tree. This is enforced by `Consistenti` and `Consistentp`.

First, `Consistenti` ensures all IDs satisfying the Presburger encoding are unique. Since we transferred these values directly to  $T$  our tree also has unique IDs.

Next, we have to ensure properties such as no node is both active and inactive. These are all directly taken care of by `Consistentp`. Thus, we are done.  $\square$

## D Additional Material for the Experiments Section

### D.1 Optimised CSS Automata Emptiness Check

The reduction presented in Section 5 proves membership in NP. However, the formula constructed is quite large even for the intersection of two relatively small selectors. Moreover, selectors generally do not assert complex properties, so for most transitions, the full power of existential Presburger arithmetic is not needed. Hence, only a small part of each formula requires complex reasoning, while the remainder of the problem is better and easily solved with direct knowledge of the automata.

In this section we present an alternative algorithm. In essence it is a backwards reachability algorithm for deciding non-emptiness of a CSS automaton. Instead of constructing a single large query that requires a non-trivial solve time, the backwards reachability algorithm only makes small queries to the SAT solver to enforce constraints that are not simply enforced by a standard automaton algorithm.

The idea is that the automaton collects constraints on the node positions required to satisfy the  $n$ th-child (sibling) constraints as it performs its backwards search. It also tracks extra information to ensure it does not get stuck in a loop, at most one node is labelled `:target`, and all `ids` are unique. Each time the automaton takes a transition labelled  $\downarrow$  it checks whether the current set of sibling constraints is satisfiable. If so, the automaton can move up to the parent node and begin with a fresh set of sibling constraints. If not, the automaton cannot execute the transition. Once the initial state has been reached, it just remains to check whether the `id` constraints are satisfiable. If they are, a witness to non-emptiness has been found.

The algorithm is a worklist algorithm, where the worklist consists of tuples of the form

$$(q, b_{\text{root}}, b_{\text{sib}}, b_{\text{targ}}, \text{Ts}, C_{\text{id}}, C_{\text{pos}}, i)$$

where

- $q$  is the state reached so far,
- $b_{\text{root}}$  is a boolean indicating whether the current node has to be the root,
- $b_{\text{sib}}$  is a boolean indicating whether the current node has to have siblings,
- $b_{\text{targ}}$  is a boolean indicating whether a node marked `:target` has been seen on the run so far,
- $\text{Ts}$  is the set of transitions seen on the current state (recall all loops are self-loops, so cycle detection can be implemented using  $\text{Ts}$ ),
- $C_{\text{id}}$  is the set of constraints on `id` attributes in the run so far,

- $C_{\text{pos}}$  is the set of constraints on node positions on the current level of the tree (that is, for the assertion of nth-child constraints),
- $i$  is the index position in the run (akin to the use of indices in the Presburger encoding).

The initial worklist contains a single element

$$(q_f, \perp, \perp, \perp, \emptyset, \emptyset, \emptyset, n)$$

where  $n$  is the number of transitions of the CSS automaton. Note, the final element of the tuple will always range between 1 and  $n$  since we decrement this counter whenever we take a new transition, and each transition may only be visited once.

In the following, we partition sets of node selector elements into sets containing pseudo-classes, attribute selectors, and positional selectors. That is, given a node selector  $\tau\Theta$  we write

- $\text{Atts}(\Theta)$  for the elements of  $\Theta$  of the form  $\theta$  or  $:\text{not}(\theta)$  where  $\theta$  is of the form  $[s|a]$  or  $[s|a\text{ op }v]$ ,
- $\text{Pos}(\Theta)$  for the elements of  $\Theta$  of the form  $\theta$  or  $:\text{not}(\theta)$  where  $\theta$  is of the form

$$\begin{aligned} &:\text{nth-child}(\alpha n + \beta), \\ &:\text{nth-last-child}(\alpha n + \beta), \\ &:\text{nth-of-type}(\alpha n + \beta), \\ &:\text{nth-last-of-type}(\alpha n + \beta) \\ &, :\text{only-child}, \text{ or} \\ &:\text{only-of-type}, \end{aligned}$$

- $\text{Pseudo}(\Theta)$  for the elements of  $\Theta$  of the form  $\theta$  or  $:\text{not}(\theta)$  where  $\theta$  is of the form

$$\begin{aligned} &:\text{link}, :\text{visited}, :\text{hover}, :\text{active}, :\text{focus}, :\text{target}, \\ &:\text{enabled}, :\text{disabled}, :\text{checked}, :\text{root}, \text{ or } :\text{empty} . \end{aligned}$$

If the worklist is empty, we terminate, and return that the automaton is empty.

If it is not empty, we take an arbitrary element

$$(q', b_{\text{root}}, b_{\text{sib}}, b_{\text{targ}}, \text{Ts}, C_{\text{id}}, C_{\text{pos}}, i)$$

and for each transition

$$t = q \xrightarrow[\sigma]{d} q'$$

with  $\sigma = \tau\Theta$  we add to the worklist

$$(q, b'_{\text{root}}, b'_{\text{sib}}, b'_{\text{targ}}, \text{Ts}', C'_{\text{id}}, C'_{\text{pos}}, i')$$

where  $i' = i-1$  is a fresh index, and when certain conditions are satisfied. We detail these conditions and the definition of the new tuple below. We begin with general conditions and definitions, then describe those specific to the value of  $d$ .

In all cases, we can only add a new tuple if

1.  $t \notin \text{Ts}$ ,
2.  $\text{AttsPres}(\tau\Theta, i')$  is satisfiable,

3.  $\text{Pseudo}(\Theta)$  is satisfiable – that is, we do not have  $p \in \Theta$  and  $:\text{not}(p) \in \Theta$  for some pseudo-class  $p$ , and, moreover, we do not have

- $:\text{link} \in \Theta$  and  $:\text{visited} \in \Theta$ , or
- $:\text{enabled} \in \Theta$  and  $:\text{disabled} \in \Theta$ , or
- $:\text{root} \in \Theta$  and  $b_{\text{sib}} = \top$ .

In all cases, we define

$$\bullet b'_{\text{root}} = \begin{cases} \top & :\text{root} \in \Theta \\ b_{\text{root}} & d = \circ \\ \perp & \text{otherwise,} \end{cases}$$

$$\bullet b'_{\text{sib}} = \begin{cases} \top & d \in \{\rightarrow, \rightarrow+\} \\ b_{\text{sib}} & d = \circ \\ \perp & \text{otherwise,} \end{cases}$$

$$\bullet b'_{\text{targ}} = b_{\text{targ}} \vee (:\text{target} \in \Theta),$$

$$\bullet \text{Ts}' = \begin{cases} \text{Ts} \cup \{t\} & q = q' \\ \emptyset & \text{otherwise,} \end{cases}$$

$$\bullet C'_{\text{id}} = C_{\text{id}} \cup C''_{\text{id}} \text{ where } C''_{\text{id}} \text{ is the set of all clauses } \text{AttsPres}_{s:\text{id}}(C, i') \text{ appearing in } \text{AttsPres}(\tau\Theta, i') \text{ for some } s \text{ and } C.$$

Next, we give the conditions and definitions dependent on  $d$ . To do so we need to define the set of positional constraints derived from  $\text{Pos}(\Theta)$ . We use a slightly different encoding to the previous section. We use a variable  $\bar{n}_i$  encoding that the node is the  $\bar{n}_i$ th child of the parent, and  $\bar{n}_i^{s:a}$  counting the number of nodes of type  $s:a$  to the left of the current node (exclusive). To encode “last of” constraints, we use the variable  $\bar{N}$  to encode the total number of siblings of the current node (inclusive), and  $\bar{N}_{s:a}$  to encode the total number of siblings of the given type (inclusive).

That is, when  $b'_{\text{root}} = \top$  let

$$\bullet C''_{\text{pos}} = \{\perp\} \text{ if } b'_{\text{sib}} = \top,$$

$$\bullet C''_{\text{pos}} = \{\perp\} \text{ if there is some } \theta \in \text{Pos}(\Theta) \text{ that is not of the form } :\text{not}(\theta') \text{ for some } \theta', \text{ and}$$

$$\bullet C''_{\text{pos}} = \{\top\} \text{ otherwise,}$$

and when  $b'_{\text{root}} = \perp$  let  $C''_{\text{pos}} =$

$$\begin{aligned}
& \{ \exists \bar{n} . \bar{n}_{i'} = \alpha \bar{n} + \beta \mid \text{nth-child}(\alpha \bar{n} + \beta) \in \Theta \} \cup \\
& \{ \exists \bar{n} . \bar{N} - \bar{n}_{i'} - 1 = \alpha \bar{n} + \beta \mid \text{nth-last-child}(\alpha \bar{n} + \beta) \in \Theta \} \cup \\
& \left\{ \begin{array}{l} \bar{s}_{i'} : \bar{e}_{i'} = s : e \Rightarrow \\ \exists \bar{n} . \bar{n}_{i'}^{s:e} = \alpha \bar{n} + \beta \end{array} \mid \begin{array}{l} s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \wedge \\ \text{nth-of-type}(\alpha \bar{n} + \beta) \in \Theta \end{array} \right\} \cup \\
& \left\{ \begin{array}{l} \bar{s}_{i'} : \bar{e}_{i'} = s : e \Rightarrow \\ \exists \bar{n} . \bar{N}_{s:e} - \bar{n}_{i'}^{s:e} - 1 = \alpha \bar{n} + \beta \end{array} \mid \begin{array}{l} s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \wedge \\ \text{nth-last-of-type}(\alpha \bar{n} + \beta) \in \Theta \end{array} \right\} \cup \\
& \{ \text{NoMatch}(\bar{n}_{i'}, \alpha, \beta) \mid \text{not}(\text{nth-child}(\alpha \bar{n} + \beta)) \in \Theta \} \cup \\
& \{ \text{NoMatch}(\bar{N} - \bar{n}_{i'} - 1, \alpha, \beta) \mid \text{not}(\text{nth-last-child}(\alpha \bar{n} + \beta)) \in \Theta \} \cup \\
& \left\{ \begin{array}{l} \bar{s}_{i'} : \bar{e}_{i'} = s : e \Rightarrow \\ \text{NoMatch}(\bar{n}_{i'}^{s:e}, \alpha, \beta) \end{array} \mid \begin{array}{l} s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \wedge \\ \text{not}(\text{nth-of-type}(\alpha \bar{n} + \beta)) \in \Theta \end{array} \right\} \cup \\
& \left\{ \begin{array}{l} \bar{s}_{i'} : \bar{e}_{i'} = s : e \Rightarrow \\ \text{NoMatch}(\bar{N}_{s:e} - \bar{n}_{i'}^{s:e} - 1, \alpha, \beta) \end{array} \mid \begin{array}{l} s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \wedge \\ \text{not}(\text{nth-last-of-type}(\alpha \bar{n} + \beta)) \in \Theta \end{array} \right\} \cup \\
& \{ \bar{N} > \bar{n}_{i'} \} \cup \left\{ \bar{n}_{i'} = \sum_{\substack{s \in \downarrow(\text{NS}) \\ e \in \downarrow(\text{ELE})}} \bar{n}_{i'}^{s:e} \right\} \cup \\
& \left\{ \begin{array}{l} \bar{s}_{i'} : \bar{e}_{i'} = s : e \Rightarrow \bar{N}_{s:e} > \bar{n}_{i'}^{s:e} \\ \bar{s}_{i'} : \bar{e}_{i'} \neq s : e \Rightarrow \bar{N}_{s:e} \geq \bar{n}_{i'}^{s:e} \end{array} \mid s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \right\} \cup \\
& \left\{ \begin{array}{l} \bar{s}_{i'} : \bar{e}_{i'} = s : e \Rightarrow \bar{N}_{s:e} > \bar{n}_{i'}^{s:e} \\ \bar{s}_{i'} : \bar{e}_{i'} \neq s : e \Rightarrow \bar{N}_{s:e} \geq \bar{n}_{i'}^{s:e} \end{array} \mid s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \right\} .
\end{aligned}$$

Then we have the following.

- When  $d = \circ$ , we set

$$\begin{aligned}
C'_{\text{pos}} &= C_{\text{pos}} \cup C''_{\text{pos}} \cup \\
& \quad \{ \bar{n}_{i'} = \bar{n}_i \} \cup \{ \bar{s}_{i'} : \bar{e}_{i'} = \bar{s}_i : \bar{e}_i \} \cup \\
& \quad \{ \bar{n}_{i'}^{s:e} = \bar{n}_i^{s:e} \mid s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \} .
\end{aligned}$$

- When  $d = \downarrow$ ,

- we require  $C_{\text{pos}}$  is satisfiable,  $\neg b_{\text{root}}$ , and  $\text{empty} \notin \Theta$ ,
- we set  $C'_{\text{pos}} = C''_{\text{pos}}$ .

- When  $d = \rightarrow$ ,

- we require  $\neg b_{\text{root}}$  and  $\text{root} \notin \Theta$ ,
- we set

$$\begin{aligned}
C'_{\text{pos}} &= C_{\text{pos}} \cup C''_{\text{pos}} \cup \{ \bar{n}_i = \bar{n}_{i'} + 1 \} \cup \\
& \quad \{ \bar{s}_{i'} : \bar{e}_{i'} = s : e \Rightarrow \bar{n}_i^{s:e} = \bar{n}_{i'}^{s:e} + 1 \mid s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \} \cup \\
& \quad \{ \bar{s}_{i'} : \bar{e}_{i'} \neq s : e \Rightarrow \bar{n}_i^{s:e} = \bar{n}_{i'}^{s:e} \mid s \in \downarrow(\text{NS}) \wedge e \in \downarrow(\text{ELE}) \} .
\end{aligned}$$

- When  $d = \rightarrow_+$

- we require  $\neg b_{\text{root}}$  and  $\text{root} \notin \Theta$ ,

– we set

$$C'_{\text{pos}} = C_{\text{pos}} \cup C''_{\text{pos}} \cup \left\{ \exists \bar{\delta}, (\bar{\delta}_{s:e})_{\substack{s \in \text{I}(\text{NS}) \\ e \in \text{I}(\text{ELE})}} \cdot \left( \bigwedge_{\substack{s \in \text{I}(\text{NS}) \\ e \in \text{I}(\text{ELE})}} \left( \begin{array}{l} \bar{n}_i = \bar{n}_{i'} + \bar{\delta} \wedge \bar{\delta} \geq 1 \wedge \\ \bar{n}_i^{s:e} = \bar{n}_i^{s:e} + \bar{\delta}_{s:e} \wedge \\ (\bar{s}_{i'} : \bar{e}_{i'} = s:e) \Rightarrow \bar{\delta}_{s:e} \geq 1 \end{array} \right) \wedge \right. \right. \\ \left. \left. \bar{\delta} = \sum_{\substack{s \in \text{I}(\text{NS}) \\ e \in \text{I}(\text{ELE})}} \bar{\delta}_{s:e} \right) \right\}.$$

Moreover, if we were able to add the tuple to the worklist, and we have

- $q = q^{\text{in}}$ ,
- $C'_{\text{pos}}$  is satisfiable, and
- The ID constraints are satisfiable,

then the algorithm terminates, reporting that the automaton is non-empty. To check the ID constraints are satisfiable, we test satisfiability of the following formula. Let  $N$  be the bound on the length of ID values, as derived in the previous section. We assert

$$\bigwedge_{\theta \in C'_{\text{id}}} \theta \wedge \bigwedge_{\substack{s \in \text{NS} \\ 1 \leq i_1, i_2 \leq n}} \bigvee_{1 \leq j \leq N} x_{i_1, j}^{s:\text{id}} \neq x_{i_2, j}^{s:\text{id}}.$$

That is, we assert all ID conditions are satisfied, and all IDs are unique.

## D.2 Sources of the CSS files used in our experiments

We have collected 72 CSS files from 41 global websites for our experiments. These websites cover the 20 most popular sites listed on Alexa [1], which are Google, YouTube, Facebook, Baidu, Wikipedia, Yahoo!, Reddit, Google India, Tencent QQ, Taobao, Amazon, Tmall, Twitter, Google Japan, Sohu, Windows Live, VK, Instagram, Sina, and 360 Safeguard. Note that we have excluded Google India and Google Japan from our collection as we found the two sites share the same CSS files with Google. We have further collected CSS files from 12 well-known websites ranked between 21 and 100 on the same list, including LinkedIn, Yahoo! Japan, Netflix, Imgur, eBay, WordPress, MSN, Bing, Tumblr, Microsoft, IMDb, and GitHub. Our examples also contain CSS files from several smaller websites, including Arch Linux, arXiv, CNN, DBLP, Google News, Londonist, The Guardian, New York Times, NetworkX, OpenStreetMap, and W3Schools. These examples were used in the testing and development of our tool.