

# WGT: Tools and algorithms for recognizing, visualizing and generating Wheeler graphs

Kuan-Hao Chao<sup>1,†,\*</sup>, Pei-Wei Chen<sup>2,†</sup>, Sanjit A. Seshia<sup>2</sup>, and Ben Langmead<sup>1,\*</sup>

<sup>1</sup>Department of Computer Science, Johns Hopkins University

<sup>2</sup>Department of Electrical Engineering and Computer Sciences,  
University of California, Berkeley

*\*corresponding author:* [kh.chao@cs.jhu.edu](mailto:kh.chao@cs.jhu.edu), [langmea@cs.jhu.edu](mailto:langmea@cs.jhu.edu)

<sup>†</sup>These authors contributed equally to this work

October 20, 2022

## Abstract

**Summary:** A Wheeler graph represents a collection of strings in a way that is particularly easy to index and query. Such a graph is a practical choice for representing a graph-shaped pangenome, and it is the foundation for current graph-based pangenome indexes. However, there are no practical tools to visualize or to check graphs that may have the Wheeler properties. Here we present *Wheelie*, an algorithm that combines a *renaming heuristic* with a Satisfiability Modulo Theory (SMT) solver to check whether a given graph has the Wheeler properties, a problem that is NP complete in general. *Wheelie* can check a variety of random and real-world graphs in far less time than any algorithm proposed to date. It can check a graph with 1,000s of nodes in seconds. We implement these algorithms together with complementary visualization tools in the WGT toolkit, available as open source software at [https://github.com/Kuanhao-Chao/Wheeler\\_Graph\\_Toolkit](https://github.com/Kuanhao-Chao/Wheeler_Graph_Toolkit).

## 1 Introduction

A Wheeler graph is a class of directed, edge-labeled graph that is particularly easy to index and query. It is a generalization of the Burrows-Wheeler-Transform-based FM Index (9), and partly forms the basis for existing pangenome alignment tools such as *vg* (11, 17).

A graph is a Wheeler Graph when its nodes can be totally ordered according to the co-lexicographical order of the sets of strings spelled out on all paths leading into the nodes. Formally: an edge-labeled, directed graph is a Wheeler graph if and only if there exists a total ordering over its nodes such that:

- 0-indegree nodes come before all other nodes in the ordering
- For all pairs of edges  $(u, v)$  and  $(u', v')$  labeled  $a$  and  $a'$  respectively:
  - $a < a' \rightarrow v < v'$
  - $a = a' \wedge u < u' \rightarrow v \leq v'$

Wheeler graphs generalize other graph- and tree-shaped structures important in genomics, including tries, De Bruijn graphs, and the reverse deterministic automata that can be constructed from multiple sequence alignments (10). These special cases of Wheeler graphs have been applied in the context of pangenome alignment tools like GCSA (16), HISAT2 (13) and VARI (14).

Despite their popularity, there are no libraries that make it easy to use Wheeler graphs, or to check if a particular graph has the requisite properties. This problem is NP-complete<sup>1</sup> and hard to approximate (12). An exponential-time algorithm was proposed in (12), but there is no implementation available.

We present WGT, an open source suite for generating, recognizing, and visualizing Wheeler graphs. WGT includes functionality for generating graphs that do or do not have the Wheeler properties. Two generators produce De Bruijn graphs and tries derived from one or more input sequences provided as FASTA. Another generator produces reverse deterministic graphs (16) from multiple sequence alignments. A third generator produces random graphs with parameterized by the desired number of nodes, edges, distinct edge labels (i.e. alphabet size), and the most number of outgoing same-label edges.

Central to WGT is the fast *Wheelie* algorithm for Wheeler graph recognition. The algorithm combines a *renaming heuristic* with two alternate solvers, both capable of reaching exact solutions to the recognition problem. One solver uses an exhaustive search over possible node permutations, and the other uses a Satisfiability Modulo Theory (SMT) solver (3). We call the overall algorithm “*Wheelie*”, while we use the names “*Wheelie-Pr*” and “*Wheelie-SMT*” for the versions that use the permutation and SMT solvers respectively. When run on a Wheeler graph, *Wheelie* also reports a node ordering for which the properties are satisfied and indexes the graph into  $O$ ,  $I$ , and  $L$  three bitarrays (10), which are useful inputs to a downstream tool for pattern matching.

Here we benchmark *Wheelie*’s solvers in comparison to each other and to the algorithm proposed by Gibney and Thankachan (12). We benchmark with a variety of input graphs, including graphs derived from real multiple alignments of DNA and protein sequences. We also use randomly-generated graphs with various configurable characteristics. Finally, we implement and demonstrate a visualizer that allows the user to picture the graph in light of the Wheeler properties.

<sup>1</sup>For special-case graphs, e.g. when at most 2 edges outgoing from a node are permitted to have the same label, the problem is tractable (1).

In the following,  $G$  denotes a directed graph,  $N$  its set of nodes, and  $E$  its set of edges, with  $n = |N|$  and  $e = |E|$ .  $\Sigma$  denotes the set of edge labels appearing on at least one edge, with  $\sigma = |\Sigma|$ .

## 2 Results

Graphs used for evaluation were generated using WGT's generator algorithms, which can produce (a) De Bruijn graphs, (b) tries, (c) a reverse deterministic graphs derived from a multiple alignments, (d) *complete* random Wheeler graphs, and (e) a *d-NFA* random Wheeler graphs. All are discussed further in Methods 3.3 and Appendix 2. For graphs that start from biological sequences, we selected 25 genes and downloaded their DNA and protein ortholog alignments in FASTA format from the Esnbml Comparative Genomics page (6). All the experiments are conducted on a 24-core, 48-thread Intel(R) Xeon(R) Gold 6248R Linux computer with 1024 GB memory, using a single thread of execution.

### 2.1 Comparing *Wheelie* with Gibney & Thanckachan

Gibney and Thankachan's recognition algorithm (12) (henceforth "G & T") works by enumerating all possible values for the  $O$ ,  $I$  and  $L$  arrays making up the Wheeler Graph structure as described by Gagie et al (10). The  $O$  bitarray is a concatenation of unary codes describing the outdegrees of each node.  $I$  is a similar bitarray that does the same for indegrees.  $L$  is a sequence of characters labeling the edges in the order they appear in the  $O$  array. Further, the inner loop of the algorithm must check if a given assignment for  $I$ ,  $O$ , and  $L$  is isomorphic to the input graph provided.

While the G & T algorithm explores an exponential-sized space, *Wheelie* explores the factorial-sized space of node permutations. While this could represent a disadvantage for *Wheelie*, we hypothesized that *Wheelie* could be made faster with the help of strategies for pruning the search space. *Wheelie* prunes its search by assigning labels to nodes according to their *rough* positions in the order, a strategy we call the "*renaming heuristic*". This allows *Wheelie* to arrive rapidly at a rough ordering that either (a) reveals a conflict that prevents the graph from having the Wheeler properties, or (b) reduces the problem size for the downstream solving algorithm. The full algorithm is described in Methods 3.1, 3.2. Here we use a version of the algorithm called *Wheelie-Pr*, which begins with the *renaming heuristic* then resolves remaining ambiguities by exhaustively searching over the remaining node permutations. Unlike G & T, there is no need to check graph isomorphisms.

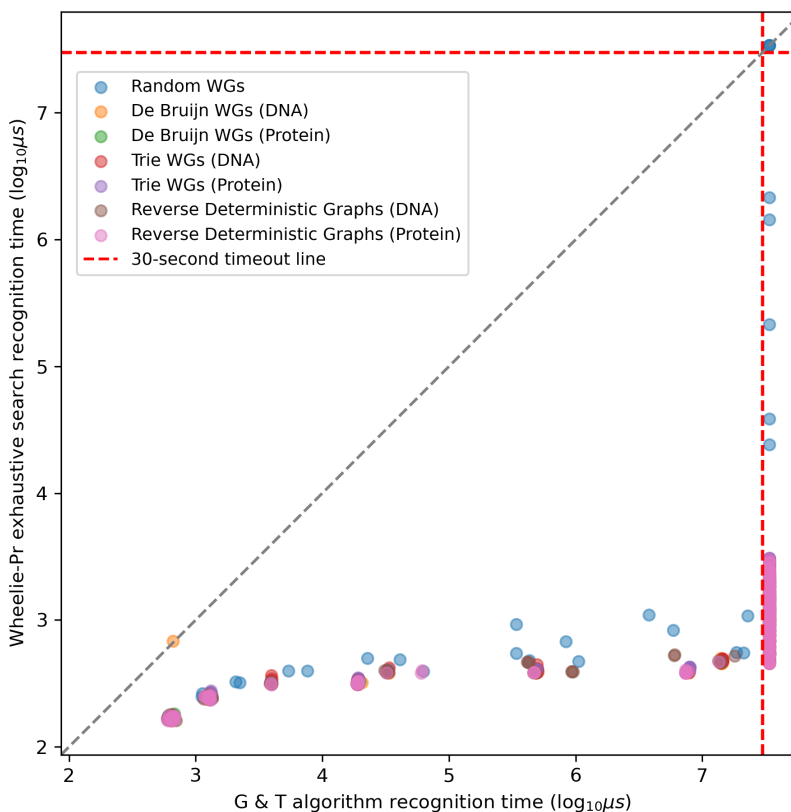
We conducted a 30-second timeout test on both algorithms using graphs generated from four generators (3.3) including both Wheeler and non-Wheeler graphs. Rather than implement G & T's entire algorithm, we implemented the enumeration of the  $I$ ,  $O$  and  $L$  arrays but omitted the graph isomorphism check in the inner loop. To compare the algorithms, we configured both to perform an exhaustive search, without the possibility of early stopping if a solution is found. This differs from *Wheelie-Pr*'s default behavior, which allows it to stop upon finding a node ordering for which the Wheeler properties are satisfied. Note that early stopping is still possible for *Wheelie-Pr* in these experiments in the event that it identifies a conflict that proves the graph is non-Wheeler.

Since G & T's exhaustive search requires exponential time, we benchmarked using small graphs. Specifically, we took a multiple alignment involving 25 genes, both their DNA and amino acid (AA) sequences, and extracted the first 4 rows of each. To reduce graph size, we also truncated the graphs with respect to the multiple-alignment columns; for De Bruijn graphs, we took columns 1 to 5 and set  $k$  to 3 and 4; for reverse deterministic graphs, we took columns 2 to 15; for tries, we took columns 2 to 10. We constructed a total of 250 DNA / 250 AA De Bruijn graphs, 350 DNA / 350 AA reverse deterministic graphs and 225 DNA / 225 AA tries. We also benchmarked the algorithms using random graphs. Specifically, we generated a series of random graphs with  $n$  set to 1 to 15,  $e$  from 1 to  $n$  and  $\sigma$  from 1 to 4, yielding 96 graphs in total.

Figure 1 shows that *Wheelie-Pr* is significantly faster, allowing it to recognize a range of Wheeler and non-Wheeler graphs. *Wheelie-Pr* runtimes generally range from 100-1,000 microseconds, with some random-graph inputs causing *Wheelie-Pr* to time out. In total only three input random graphs caused *Wheelie-Pr* to time out whereas 1,748 input graphs caused G & T to time out.

### 2.2 Visualizing and characterizing challenging graphs

We sought to better understand which graphs require the most time for recognition. We selected a De Bruijn graph with edges being  $k$ -mers and nodes being  $k-1$ -mers where  $k = 4$  from the Figure1 benchmarks. This graph was derived from the first four rows of the multiple alignment of STAU2 DNA orthologues with



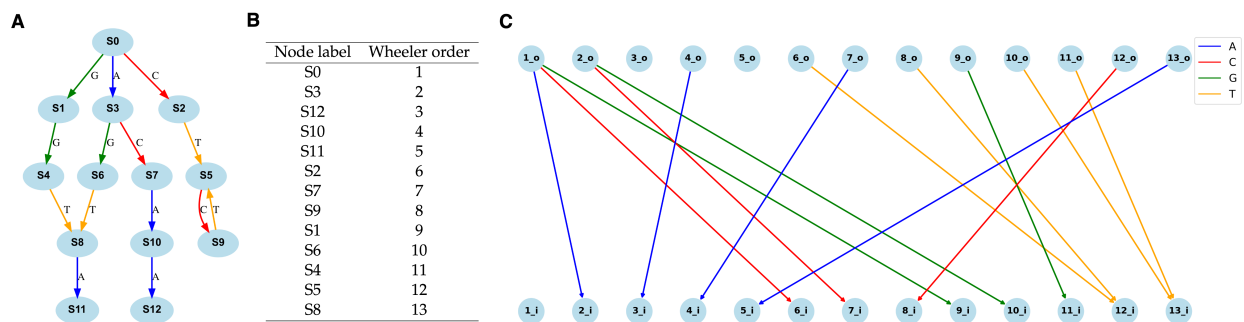
**Figure 1:** Recognition time comparison between `Wheelie-Pr` exhaustive search and the G & T algorithm using (1) De Bruijn graphs, (2) tries and (3) reverse deterministic graphs generated from DNA and protein alignments, and (4) random graphs generated with given  $n$ ,  $e$  and  $\sigma$ . `Wheelie-Pr` is on the y-axis, and G & T algorithm is on the x-axis. Both axes are in  $\log_{10}$  microsecond scale. Each dot represent a graph, and dots beyond the red lines denote inputs for which the tool timed out after 30 seconds.

sequence length 4. We first visualized it by Graphviz online visualizer (8) (Figure 2A). We ran `Wheelie-Pr` to find an ordering for which the Wheeler properties hold (Figure 2B). Finally, we visualized the graph using WGT’s Python-based visualizer, which draws the ordered nodes in two replicas, with outgoing edges leaving one replica and entering the other. This is a useful way to visualize and validate recognition results: for a valid Wheeler ordering, nodes with no incoming edges will appear leftmost, nodes with incoming edges of the smallest character will come next, nodes with incoming edges of the next-smallest character next, etc. Additionally, no two same-color edges will cross each other. In this way, the diagram – used previously by Boucher et al (5) – makes it visually obvious when an ordering has yielded the Wheeler properties.

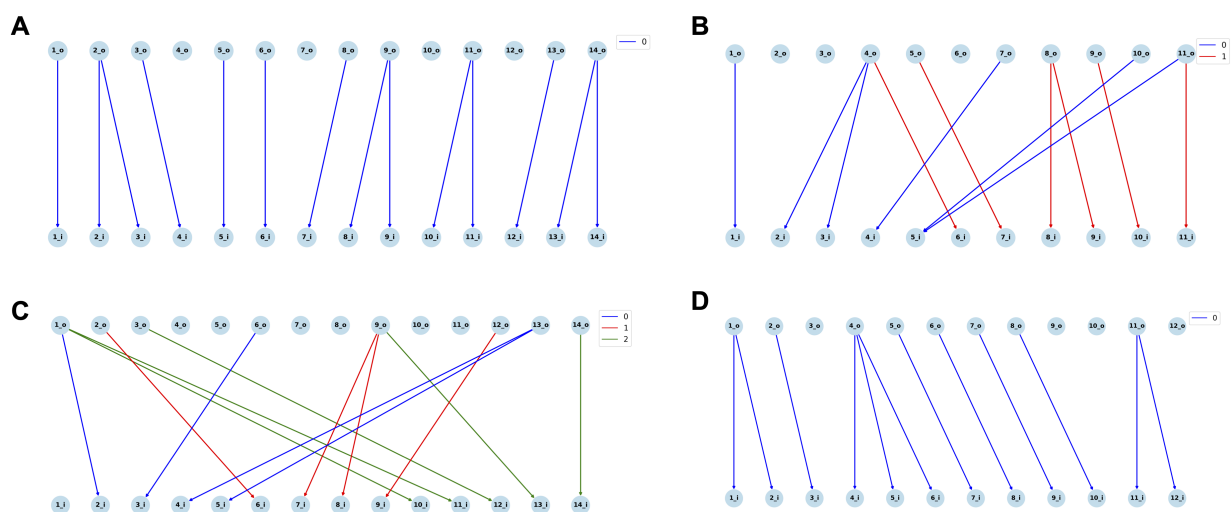
After investigating with these tools, we found that the graphs requiring the most recognition time tended to have nodes with many outgoing same-label edges. Following Alanko et al (1), we use the term  $d$ -NFA to describe a Wheeler Graph where all nodes have  $\leq d$  outgoing same-label edges, and at least one node has exactly  $d$  outgoing same-label edges. The De Bruijn graph shown in Figure 2 is a 1-NFA. Figure 3 A-C are 2-NFAs with  $\sigma$  equal to 1, 2 and 3 respectively. Figure 3 D is a 3-NFA with  $\sigma = 1$ . Note that the case where  $d \geq 5$  is the one that has been proven to be NP-complete (12). De Bruijn graphs and tries are 1-NFAs.

### 2.3 Recognizing challenging graphs with `Wheelie-SMT`

Motivated by previous work that showed how boolean satisfiability formulations can solve special cases of the recognition problem (1), we hypothesized that Satisfiability Modulo Theories (SMT) solvers (3) could be used to solve all or part of the Wheeler-graph recognition problem. SMT has found many uses in artificial intelligence and formal methods for hardware and software development. As a generalization of the Boolean Satisfiability (SAT) (4), SMT allows us to encode the Wheeler graph properties in a fairly straight-



**Figure 2:** (A) A  $k = 4$  De Bruijn graph outputted from WGT’s De Bruijn graph generator. It is the visualization from Graphviz online visualizer. (B) The recognition result showing the Wheeler ordering outputted from *Wheelie-Pr*. (C) The output from WGT’s visualizer. Nodes are duplicated into two rows ordered in Wheeler ordering.



**Figure 3:** Examples of (A) 2-NFA with 1 label (B) 2-NFA with 2 labels (C) 2-NFA with 3 labels (D) 3-NFA with 1 label, from outlier random graphs (blue dots) in Figure 1.

forward way, building from the propositional logic formulas in the definition.

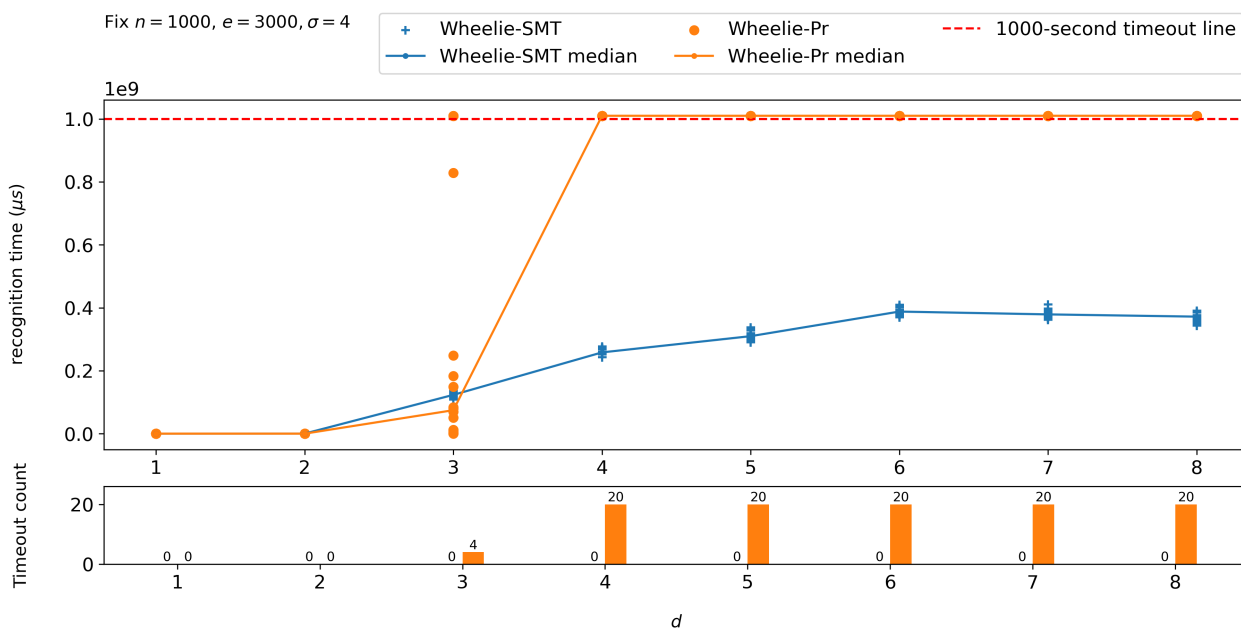
We conducted two series of 1,000-second timeout tests using graphs generated from the random generator comparing (1) *Wheelie-Pr* (*renaming heuristic* plus permutation) versus (2) *Wheelie-SMT* (*renaming heuristic* plus SMT) on different types of  $d$ -NFA (2.3.1) and various sizes of random graphs (2.3.2).

### 2.3.1 Recognizing $d$ -NFAs

We fixed  $n = 1000$ ,  $e = 3000$  and  $\sigma = 4$  and randomly generated  $d$ -NFAs with  $d$  from 1 to 8 and each group with 20 graphs. Figure 4 shows that both solvers can solve graphs swiftly when  $d$  is 1 and 2; as  $d$  grows beyond 2, all tools require much more time, demonstrating that  $d$  impacts the hardness of recognition problem in practice. *Wheelie-SMT* outperforms *Wheelie-Pr* and avoids any timeouts; *Wheelie-Pr* has some timeouts starting at  $d = 3$  (4 out of 20 graphs), and consistently times out when  $d \geq 4$ .

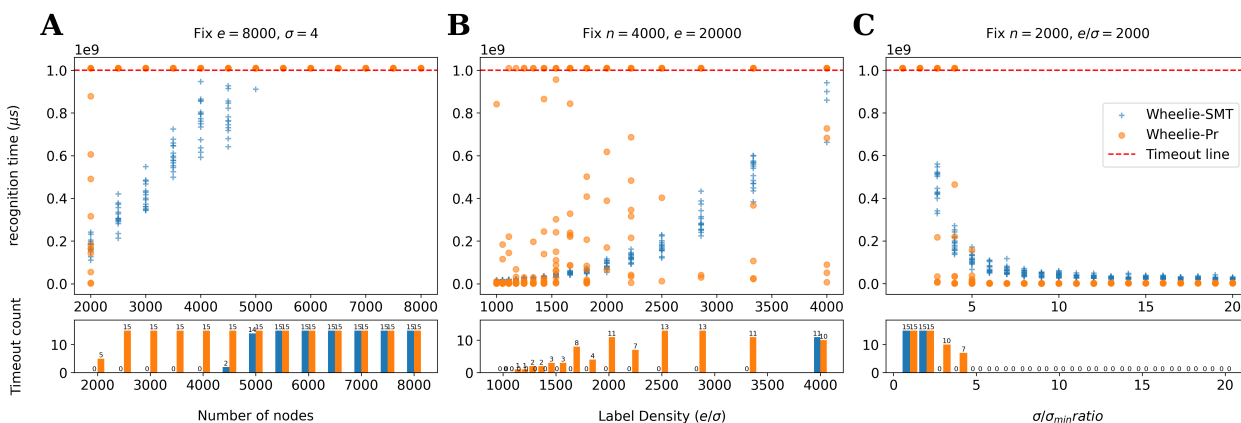
Further, we observed that when  $d \geq 6$ , the median curve for *Wheelie-SMT* plateaus. This is because  $n$  and  $e$  are too small for the  $d$ -NFA generator to produce uniformly distributed  $d$ -NFAs under the given parameters. More precisely speaking, the hardness of the recognition problem is a function of the distribution of nodes having  $d - 1$ ,  $d - 2$ , ..., 1 out-going edges with the same labels. As an example, take a  $d$ -NFA  $G$  that has one node with  $d$  same-label out-going edges, and the rest of the nodes having at most one outgoing same-label edge. Recognizing  $G$  is not harder than recognizing a uniformly-distributed  $d - 1$ -NFA. In short, we observed that higher  $d$ s generally led to a harder recognition problem, but the true level of hardness was

also a function of  $n$ ,  $e$  and  $\sigma$ .



**Figure 4:** Recognition time for *Wheelie-SMT* and *Wheelie-Pr* as a function of the  $d$  parameter of the  $d$ -NFA. Upper panel plots recognition time versus  $d$  and includes a line connecting the medians. 20 graphs were tested for each  $d$ . The bottom bar chart shows the number of timeouts.

### 2.3.2 Recognizing random Wheeler graphs



**Figure 5:** Recognition time comparison between *Wheelie-SMT* and *Wheelie-Pr* under various random Wheeler graphs. Three experiments were conducted. (A) Experiment 1: fixing  $e = 8,000$ ,  $\sigma = 4$  and scaling up the graph size ( $n$  from 2,000 to 8,000). (B) Experiment 2: fixing  $n = 4,000$ ,  $e = 20,000$  and scaling up the label density ( $e/\sigma$  from 1,000 to 4,000). (C) Experiment 3: Fixing both graph size ( $n$ ) and label density ( $e/\sigma$ ), and scaling up both  $e$  and  $\sigma$  ( $\sigma/\sigma_{min}$  from 1 to 20). The upper-panel plots show the recognition time versus the scale up parameter in microsecond scale. Each dot represents a graph. Dots beyond the red dashed line means timeouts. Plots in the lower panel are the timeout count bar charts.

We defined “graph size” as  $n$  and “label density” as  $e/\sigma$ . We then benchmarked various sizes of random graphs while varying these parameters.

First we fixed the number of edges ( $e = 8,000$ ) and the number of labels ( $\sigma = 4$ ), and scaled up the “graph size” ( $n$  from 2,000 to 8,000). Figure 5A shows that as  $n$  grows, *Wheelie-SMT* outperforms



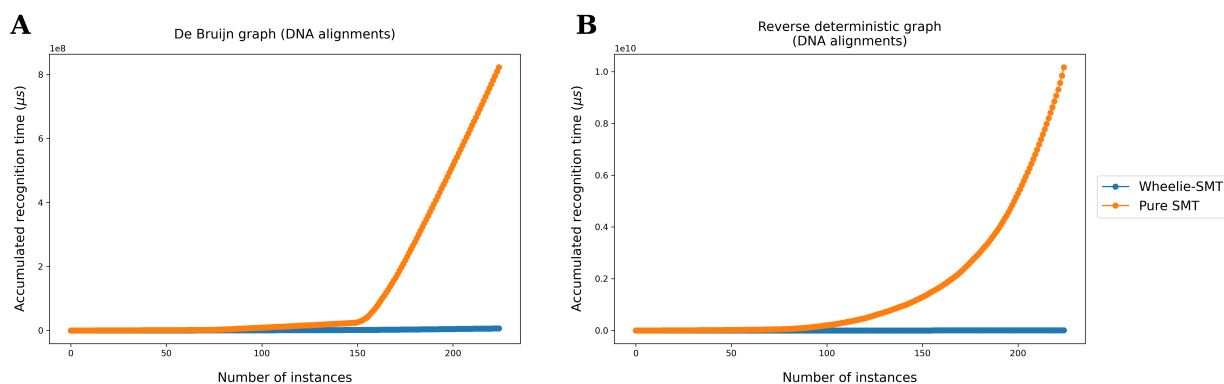
Wheelie-Pr significantly. Wheelie-Pr starts to time out in some cases when  $n = 2,000$ , and most cases when  $n \geq 2,500$ . In contrast, Wheelie-SMT can solve all cases with  $n$  up to 4,000, and most cases when  $n = 4,500$ .

We then fixed the “graph size” ( $n = 4,000$ ) and number of edges ( $e = 20,000$ ) and varied the “label density” ( $e/\sigma$  from 1,000 to 4,000). Figure 5B shows that as the label density increases, the graphs take more time to solve. Comparing Wheelie-Pr and Wheelie-SMT, we can see that there are more timeout cases in Wheelie-Pr from 1,200 to 4,000 (most are timeouts when  $e/\sigma \geq 2,500$ ) whereas the timeout cases only occur in Wheelie-SMT when label density is 4,000.

In a third experiment, we fixed the graph size ( $n = 2,000$ ) and varied the number of edges ( $e$ ) and labels ( $\sigma$ ) while fixing the label density ratio ( $e/\sigma = 2,000$ ). Figure 5C shows that as more edges and labels are added, the recognition problem becomes easier. In short, this is because adding more constraints to  $G$  breaks more of the ties that would otherwise obstruct Wheelie’s *renaming heuristic*. Comparing Wheelie-Pr to Wheelie-SMT, Figure 5C shows that the solvers perform similarly, with Wheelie-Pr performing slightly better when  $\sigma/\sigma_{min}$  ratio gets larger ( $\geq 5$ ). These are likely cases where the graph is sufficiently easy to recognize that the overhead of setting up the SMT setup problem becomes harmful. When  $\sigma/\sigma_{min}$  gets smaller (3 and 4), Wheelie-SMT is able to solve all 15 cases, whereas Wheelie-Pr’s times out for about half the cases.

## 2.4 Benchmarking Wheelie-SMT alone

To isolate the effect of the Wheelie *renaming heuristic*, we conducted a 30-second timeout test with 60 seconds timeout penalties on (1) Wheelie-SMT (*renaming heuristic* plus SMT) and (2) a pure SMT solver starting from scratch, without the constraints it would otherwise receive from the *renaming heuristic*. We benchmarked these using two generators from DNA alignments: (1) De Bruijn graphs generated with options  $-k$  from 5 to 8,  $-l$  from 100 to 2,000 and  $-a$  from 6 to 10, 225 graphs in total and (2) reverse deterministic graphs generated with options  $-l$  from 100 to 500 and  $-a$  from 4 to 6, in total 225 graphs.



**Figure 6:** The cactus survival plots of (A) De Bruijn graphs and (B) reverse deterministic graphs generated from DNA alignments using WGT’s generators. They show the aggregated time comparison between Wheelie-SMT and pure SMT without the constraints from the *renaming heuristic*.

Figure 6 shows cactus plots on De Bruijn graphs and reverse deterministic graphs. Cactus plot is an aggregated sorted time plots widely used in solver competitions. It shows how many problems a solver can solve in a limited time period. In Figure 6A, Wheelie-SMT solved the whole De Bruijn graph set in around 6.5 seconds whereas the pure SMT approach solved it in around 820 seconds. For reverse deterministic graphs (Figure 6B), Wheelie-SMT solved the whole set in less than 9 seconds whereas the pure SMT approach solved it in around 10,170 seconds.

We concluded that the *renaming heuristic* is a crucial step, since it greatly narrows the space of possible node ordering that must be resolved by the SMT solver. Wheelie-SMT can solve graphs several orders of magnitude larger than a pure SMT approach.

## 3 Methods

### 3.1 *Wheelie* and the *renaming heuristic*

*Wheelie* explores the space of possible node orderings until arriving either at a conflict (e.g. a node with distinctly labeled incoming edges) or an ordering for which the Wheeler properties hold. While this is a large ( $n!$ -sized) search space, *Wheelie* prunes the space by assigning labels to nodes according to their *rough* position in the overall order. Initially, a rough ordering is determined according to the labels of the immediate incoming edges for each node, following the Wheeler requirement that  $a \prec a' \rightarrow v < v'$  for all edge pairs. This rough ordering is refined over the course of a procedure that iterates either until the rough ordering becomes total ordering, or until the rough ordering stabilizes. In the latter case, the remaining ambiguities are resolved by a non-heuristic solver. This procedure is detailed in Algorithm 1 and illustrated in Figure 7.

As the *renaming heuristic* iterates, it repeatedly visits the nodes in groupings according to the label of their incoming edge(s). For each of these groupings, it sorts the edges by sources and destinations in every label group, requiring  $O(\prod_{g \in \Sigma} e(g) \log_2 e(g))$  time, where  $e(g)$  is the number of edges labelled as  $g$ . We observed that many non-Wheeler graphs can be recognized as such directly by the *renaming heuristic*, without requiring a downstream solver.

At each iteration, the algorithm gathers a list of sorted unique temporary orders of nodes that go into it, which we term the “in-node list.” By the Wheeler graph property that requires all edge pairs to satisfy  $a = a' \wedge u < u' \rightarrow v \leq v'$ , we can find rough orders by sorting the nodes by their in-node lists. Once this has been done for each node group, we reach the end of the current iteration and we check if the rough order changed since the previous iteration. If not, then we say the algorithm has converged and forward any remaining ambiguities to the downstream solver as necessary.

### 3.2 Satisfiability Modulo Theories (SMT) solver

Motivated by the use of boolean satisfiability formulations to solve special cases of the recognition problem (1), we hypothesized that Satisfiability Modulo Theory (SMT) solvers (3) could be used to solve all or part of the Wheeler-graph recognition problem. SMT has found many uses in artificial intelligence and formal methods for hardware and software development. As a generalization of the Boolean Satisfiability (SAT) (4), SMT allows us to encode the Wheeler graph properties in a fairly straightforward way, building from the propositional logic formulas in the definition.

An SMT problem decides the satisfiability of a first-order formula with respect to one or more background theories. A *formula* is a set of atoms connected by Boolean connectives ( $\wedge, \vee, \neg$ ), where an *atom* is a predicate which evaluates to True or False given an assignment to the variables. A *literal* is either an atom or its negation. A *theory* gives special meanings, known as *interpretations*, to functions and predicate symbols within the theory. In this paper, we consider only *the theory of Integer Difference Logic (IDL)*, which requires atoms to be of the form  $x_i - x_j \leq c$ , where  $x_i$  and  $x_j$  are integer variables,  $c$  an integer constant, “ $-$ ” the integer subtraction function, and “ $\leq$ ” the usual binary ordering predicate. A theory solver decides the satisfiability of a conjunction of *literals*. In particular, an IDL theory solver can be implemented as the Bellman-Ford algorithm which runs in polynomial time. Incorporating a SAT solver and a theory solver, an SMT solver takes in a formula and outputs an assignment to the variables if the formula is satisfiable or otherwise reports unsatisfiability.

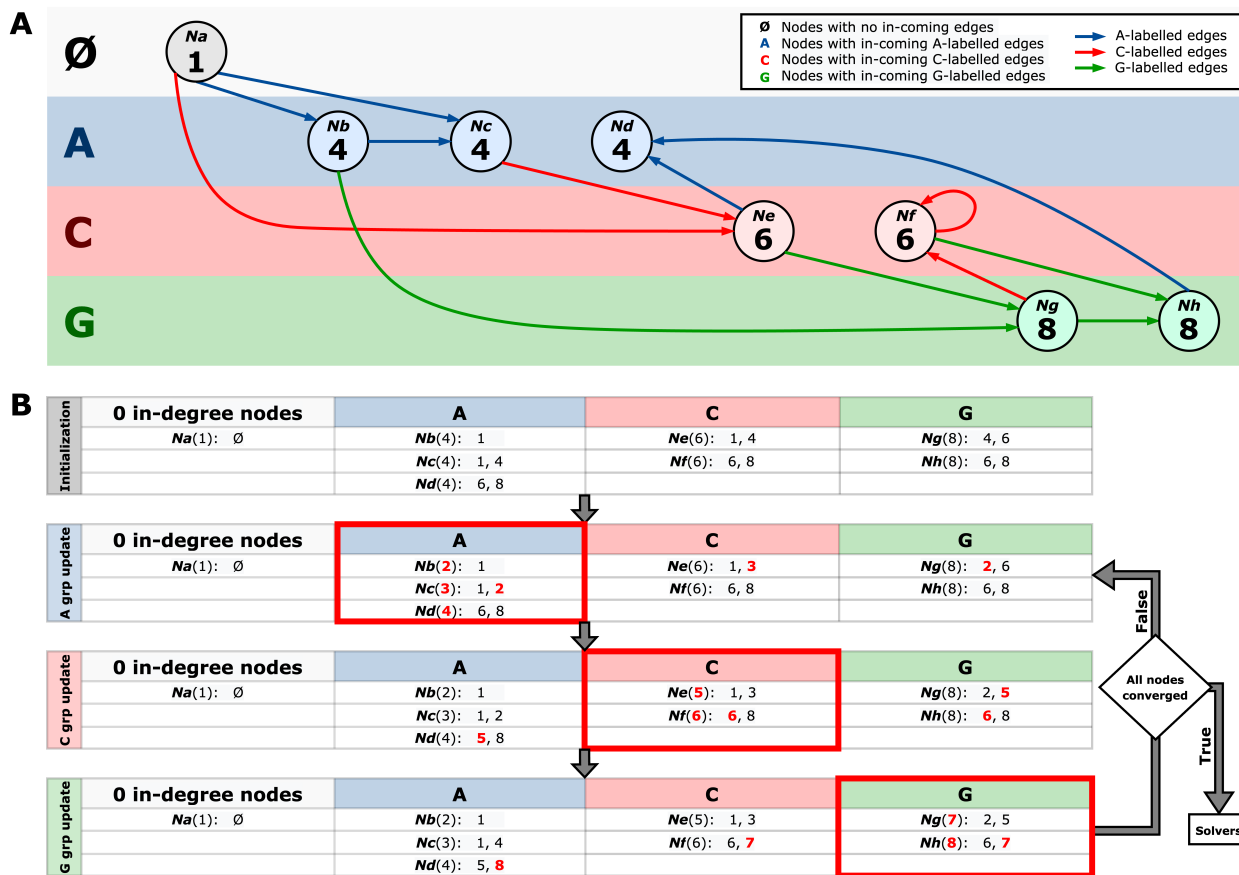
We observed that SMT is a natural way of encoding the Wheeler graph recognition problem. Firstly, for each node a variable is created representing the ordering of the node. Recalling the constraints, for all pairs of edges  $(u, v)$  and  $(u', v')$  labeled  $a$  and  $a'$  respectively:

$$a \prec a' \rightarrow v < v', \quad (1)$$

$$a = a' \wedge u < u' \rightarrow v \leq v'. \quad (2)$$

By indexing the known labels lexicographically, we obtained an SMT formula containing constraints (1) and (2), in which all atoms are of the form  $v \leq v'$  satisfying the IDL requirement. Note that the strict inequality  $v < v'$  can be rewritten using the non-strict one as  $v - v' \leq -1$ . Besides constraints (1) and (2), we





**Figure 7:** Illustration of the *renaming heuristic*. (A) an 8-node graph with nodes divided into four groups according to in-coming edge label (with  $\emptyset$  representing 0-indegree nodes). (B) presents the workflow of the *renaming heuristic*. The first table in (B) shows the initialized in-node lists for eight nodes. After initialization, the algorithm sorts and relabels nodes in each group until convergence. Then, it passes the range information to either *Wheelie-Pr* or *Wheelie-SMT*.

also enforced all-different constraints and range constraints for all nodes, which have the form  $1 \leq v \leq n$ . The node orderings can be obtained from the satisfying assignment by solving the SMT formula iff it is satisfiable. In *Wheelie-SMT*, we used the Z3 Theorem Prover (7) to encode the Wheeler Graph constraints.

As the number of constraints and variables are the main factors affecting runtime, simplifying the problem or providing additional information can improve performance. We noticed that the problem can be simplified using the rough order obtained from the *renaming heuristic* described in Section 3.1. By adding the range information to the formula, notice that constraint (1) can be removed from the formula. Moreover, the all-different constraints of nodes from different groups can also be removed. The rationale is that if the graph was not reported non-Wheeler during the *renaming heuristic*, the range constraints provided by the procedure must automatically imply constraint (1). The benefits are twofold: not only the number of constraints is reduced, the search space is also significantly pruned.

### 3.3 WGT's graph generating algorithms

We implemented five generators in Python scripts to produce tries, reverse deterministic graphs (16), De Bruijn graphs and *complete* and *d-NFA* random Wheeler graphs. The first three generators take either DNA or protein multiple sequence alignments and produce the corresponding graph structures. As for the two random generators, users can produce a Wheeler graph given its  $n$ ,  $e$ , and  $\sigma$ , and *d-NFA* random generator can further take  $d$ , which controls the most number of edges coming out from a node with the same label (*d-NFA*) as user input.

---

**Algorithm 1** Wheeler graph Recognition Algorithm, *Wheelie*

---

**Require:**

**Input:** Graph  $G$  as DOT file  
**Input:** Solver  $S$  input from `-s` or `--solver` tag.

- 1: Find all 0-indegree nodes,  $Roots[]$
- 2: Group edges by their labels in a hashMap,  $label\_2\_edges$  (key: label, value: list of edges)
- 3:  $accum\_order \leftarrow Roots.size()$  ▷ Relabel nodes with the largest possible order
- 4: **for** each  $root \in Roots$  **do**
- 5:      $relabel\_node(root, accum\_label)$
- 6: **end for**
- 7: **for** each  $[label, edges] \in label\_2\_edges$  **do**
- 8:      $accum\_order \leftarrow accum\_order + unique(edges.destination\_nodes).size()$
- 9:     **for** each  $edge \in edges$  **do**
- 10:          $relabel\_node(edge.destination\_node, accum\_order)$
- 11:     **end for**
- 12: **end for**
- 13:  $converged \leftarrow False$  ▷ renaming heuristic
- 14:  $ranges[]$
- 15: **while** not  $converged$  **do**
- 16:      $ranges.clear()$
- 17:      $accum\_order \leftarrow 0$
- 18:      $innode\_list[]$
- 19:     **for** each  $[label, edges] \in label\_2\_edges$  **do**
- 20:          $nodes \leftarrow unique(edges.destination\_nodes)$
- 21:         **for** each  $node \in nodes$  **do**
- 22:              $innode\_list.Add(get\_innodelist(node))$  ▷ get\\_innodelist: lists distinct predecessor nodes in order by label
- 23:         **end for**
- 24:          $indices = sort\_node\_by\_innodelist(nodes, innode\_list)$
- 25:          $relabel\_node(nodes, innode\_list, indices, accum\_order)$  ▷ Assign order\\_range to prev\\_order\\_range for each node
- 26:          $accum\_order \leftarrow accum\_order + nodes.size()$
- 27:     **end for**
- 28:      $converged \leftarrow True$
- 29:     **for** each  $node \in nodes$  **do**
- 30:          $ranges.Add(node.order\_range)$
- 31:         **if**  $node.order\_range \neq node.prev\_order\_range$  **then**
- 32:              $converged \leftarrow False$
- 33:              $Break$
- 34:         **end if**
- 35:     **end for**
- 36: **end while**
- 37: **if**  $unique(ranges).size() = ranges.size()$  **then** ▷ Solved by renaming heuristic
- 38:      $G$  is a Wheeler graph
- 39: **else**
- 40:     **if**  $S == "Permutation"$  **then**
- 41:          $Permutation(G, ranges)$  ▷ Use *Wheelie-Pr* solver to resolve multi-node groups
- 42:     **else if**  $S == "SMT"$  **then**
- 43:          $SMT(G, ranges)$  ▷ Use *Wheelie-SMT* solver to resolve multi-node groups
- 44:     **end if**
- 45: **end if**

---

Tries and De Bruijn graphs are Wheeler graphs by definition. These generators start by removing gap placeholders from the multiple alignments. The trie generator iterates through the prefixes of each sequence in the multiple alignment, inserts characters into the trie and creates a new node at the end of a path if a prefix cannot be traversed from the source. Edges are labelled according to the label of the parent. The De Bruijn graph generator constructs a distinct  $k - 1$ -mer dictionary from the sequences. It connects edges between adjacent two nodes and label the edge with the first character in the  $k - 1$ -mer of the child node. Reverse deterministic graphs are usually invalid Wheeler graphs but might be valid when the graphs are small, and once violations occur, adding more nodes and edges cannot turn them back to Wheeler graphs.

The reverse deterministic graph generator iterates through columns of a multiple sequence alignments from right to left. At a column  $i$ , it creates distinct nodes for the characters found there, connecting them to the current node with the node of the previous ungapped character with the direction pointing to the end of the alignments and the label of the previous ungapped character. This follows the procedure described in the GCSA study (16).

Last, three generators initializes the names of nodes with the breadth first search orders and outputs the constructed graph in DOT format.

For the two random generators, *complete* Wheeler graph and *d-NFA* Wheeler graph generators, we explain our design concepts in Appendix section2.

## 4 Discussion

We demonstrated that `Wheelie-SMT` is the fastest and most robust algorithm available for the Wheeler Graph recognition problem. We showed this across a variety of graph types, including large graphs (thousands of nodes and edges) and challenging graphs, such as those that are *d*-NFAs with values of  $d$  up to 8. We also demonstrated WGT's facilities for visualizing and understanding these graphs.

When `Wheelie` determines that a graph is is Wheeler Graph, it is able to report a node ordering that can then be used to index the graph. In the future, it will be important to extend `Wheelie` to report other useful information, including when the graph is not a Wheeler Graph. For instance, when `Wheelie` encounters a conflict that proves the graph to be non-Wheeler, `Wheelie` could supply the user with an explanation for why the graph cannot be Wheeler. Such an explanation could also allow `Wheelie` to suggest modifications to the graph that would make it a Wheeler Graph, without changing which strings it encodes. A trivial example would be a node with two incoming edges having two distinct labels. This violates the Wheeler Graph properties, but also suggests a potential solution: the node could be duplicated, with outgoing edges also duplicated. The initial inbound edges could be redrawn to point to the distinct duplicates, possibly restoring the Wheeler properties. A more general approach for understanding Wheeler violations could work by extracting conflicting sets of clauses from the SMT algorithm, and converting them into a human-understandable or other actionable form.

It may also be possible to encode the *renaming heuristic* as a set of clauses in the SMT solver, potentially allowing the entire algorithm to execute within the SMT solver. Finally, as different SMT solvers such as CVC5 (2) or Z3 (7) adopt different heuristics, they could potentially be substituted into WGT, or combined for increased efficiency (15).

## 5 Acknowledgments

We thank Markus J. Sommer for proposing the name `Wheelie` for our Wheeler graph recognition algorithm. We thank Travis Gagie and Christina Boucher for helpful comments.

## 6 Funding

This research was supported in part by the U.S. National Institutes of Health under grant R01-HG006677 and grants R35GM139602 and R01HG011392 to BL. This work was also supported by the U.S. National Science Foundation under grant DBI-1759518, and by a Berkeley Fellowship.

## 7 Authors' contributions

KC, PC and BL designed the method. KC and PC wrote the software and performed the experiments. KC, PC, SAS and BL wrote the manuscript. All authors read and approved the final manuscript.

## References

- [1] Alanko, J., D’Agostino, G., Policriti, A., and Prezza, N. (2020). Regular languages meet prefix sorting. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 911–930. SIAM.
- [2] Barbosa, H., Barrett, C. W., Brain, M., Kremer, G., Lachnitt, H., Mann, M., Mohamed, A., Mohamed, M., Niemetz, A., Nötzli, A., Ozdemir, A., Preiner, M., Reynolds, A., Sheng, Y., Tinelli, C., and Zohar, Y. (2022). cvc5: A versatile and industrial-strength SMT solver. In *28th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, TACAS 2022*, volume 13243 of *Lecture Notes in Computer Science*, pages 415–442. Springer.
- [3] Barrett, C., Sebastiani, R., Seshia, S. A., and Tinelli, C. (2009). Satisfiability modulo theories. In A. Biere, H. van Maaren, and T. Walsh, editors, *Handbook of Satisfiability*, chapter 26, pages 825–885. IOS Press.
- [4] Biere, A., Heule, M., and van Maaren, H. (2009). *Handbook of satisfiability*, volume 185. IOS press.
- [5] Boucher, C., Gagie, T., Kuhnle, A., Langmead, B., Manzini, G., and Mun, T. (2019). Prefix-free parsing for building big BWTs. *Algorithms Mol Biol*, **14**, 13.
- [6] Cunningham, F., Allen, J. E., Allen, J., Alvarez-Jarreta, J., Amode, M. R., Armean, I. M., Austine-Orimoloye, O., Azov, A. G., Barnes, I., Bennett, R., et al. (2022). Ensembl 2022. *Nucleic acids research*, **50**(D1), D988–D995.
- [7] de Moura, L. M. and Bjørner, N. S. (2008). Z3: an efficient SMT solver. In *Proc. 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems TACAS 2008*, volume 4963 of *Lecture Notes in Computer Science*, pages 337–340. Springer.
- [8] Ellson, J., Gansner, E., Koutsofios, L., North, S. C., and Woodhull, G. (2001). Graphviz—open source graph drawing tools. In *International Symposium on Graph Drawing*, pages 483–484. Springer.
- [9] Ferragina, P. and Manzini, G. (2000). Opportunistic data structures with applications. In *Proceedings 41st annual symposium on foundations of computer science*, pages 390–398. IEEE.
- [10] Gagie, T., Manzini, G., and Sirén, J. (2017). Wheeler graphs: A framework for bwt-based data structures. *Theoretical computer science*, **698**, 67–78.
- [11] Garrison, E., Sirén, J., Novak, A. M., Hickey, G., Eizenga, J. M., Dawson, E. T., Jones, W., Garg, S., Markello, C., Lin, M. F., Paten, B., and Durbin, R. (2018). Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nat Biotechnol*, **36**(9), 875–879.
- [12] Gibney, D. and Thankachan, S. V. (2019). On the hardness and inapproximability of recognizing wheeler graphs. *arXiv preprint arXiv:1902.01960*.
- [13] Kim, D., Langmead, B., and Salzberg, S. L. (2015). Hisat: a fast spliced aligner with low memory requirements. *Nature methods*, **12**(4), 357–360.
- [14] Muggli, M. D., Bowe, A., Noyes, N. R., Morley, P. S., Belk, K. E., Raymond, R., Gagie, T., Puglisi, S. J., and Boucher, C. (2017). Succinct colored de Bruijn graphs. *Bioinformatics*, **33**(20), 3181–3187.
- [15] Pimpalkhare, N., Mora, F., Polgreen, E., and Seshia, S. A. (2021). MedleySolver: Online SMT algorithm selection. In *24th International Conference on Theory and Applications of Satisfiability Testing (SAT)*, volume 12831 of *Lecture Notes in Computer Science*, pages 453–470. Springer.
- [16] Sirén, J., Välimäki, N., and Mäkinen, V. (2014). Indexing graphs for path queries with applications in genome research. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, **11**(2), 375–388.
- [17] Sirén, J., Monlong, J., Chang, X., Novak, A. M., Eizenga, J. M., Markello, C., Sibbesen, J. A., Hickey, G., Chang, P. C., Carroll, A., Gupta, N., Gabriel, S., Blackwell, T. W., Ratan, A., Taylor, K. D., Rich, S. S., Rotter, J. L., Haussler, D., Garrison, E., and Paten, B. (2021). Pangenomics enables genotyping of known structural variants in 5202 diverse genomes. *Science*, **374**(6574), abg8871.

# Appendix of

## “WGT: Tools and algorithms for recognizing, visualizing and generating Wheeler graphs”

Kuan-Hao Chao<sup>1,†,\*</sup>, Pei-Wei Chen<sup>2,†</sup>, Sanjit A. Seshia<sup>2</sup>, and Ben Langmead<sup>1,\*</sup>

<sup>1</sup>Department of Computer Science, Johns Hopkins University

<sup>2</sup>Department of Electrical Engineering and Computer Sciences,  
University of California, Berkeley

\**corresponding author*: [kh.chao@cs.jhu.edu](mailto:kh.chao@cs.jhu.edu), [langmea@cs.jhu.edu](mailto:langmea@cs.jhu.edu)

<sup>†</sup>These authors contributed equally to this work

October 20, 2022

## 1 Wheelie’s permutation based approach

While the G & T algorithm explores an exponential-sized space of possible array assignments, *Wheelie* explores a factorial-sized space of node permutations. This may or may not lead to a larger search space for *Wheelie*, depending on the graph’s properties. To be specific, the G & T’s algorithm may have to consider all  $2^{2(e+n)+e \log(\sigma)}$  assignments for  $I$ ,  $O$  and  $L$ . Our approach might need to consider  $n!$  node permutations in the worst case. We sought a rough comparison between the approaches in light of the fact that G & T’s space depends not only on  $n$  but also on  $e$  and  $\sigma$ . Below Derivation1, we fixed  $n$  for both, defining a new variable  $C$  as  $e(2+\log \sigma)$ . We then found some values for  $C$  that equalize the algorithms’ search space size under various values for of  $ns$  (Table 1). For instance, when  $n = 100$ ,  $C$  can be at most 324 in order for G & T’s algorithm has an equal or smaller search space than *Wheelie-Pr*, which is a strict threshold, and furthermore, this comparison is done with *Wheelie-Pr* skipping the *renaming heuristic*, which in reality makes *Wheelie-Pr* superiorly faster (Results2.4).

To gain a further advantage over the G & T algorithm, *Wheelie* further strives to prune the search space, using a renaming heuristic, an SMT solver, or both, as detailed in Methods3.

		n	C threshold
$2^{2(n+e)+e \times \log \sigma} = n!$	(1)	10	1.79
$2^{2n+e \times (2+\log \sigma)} = n!$	(2)	20	21.08
		30	47.71
		40	79.16
		60	152.13
Let $C := e(2 + \log \sigma)$		80	234.83
$2n + C = \log_2 n!$	(3)	100	324.76
		150	572.86
$C = \sum_{x=1}^n \log_2 x - 2n$	(4)	200	845.38

Derivation 1: The relationship between  $C$  and  $n$

**Table 1:** Threshold values of  $C$  as a function of  $n$  such that values greater than the threshold cause the permutation-based approach to have a smaller search space compared to the G & T approach.

## 2 Random generators

We implemented two random generators, a *complete* Wheeler graph generator and a *d-NFA* Wheeler graph generator. We first fix the ordering of nodes and then try to select edges such that both user-specified constraints and Wheeler graph properties are satisfied. Let  $N_i$  be the nodes with incoming edges labeled  $i$  and  $E_i$  be the edges labeled  $i$  where  $i = 1, 2, \dots, \sigma$ , and also let  $n_i = |N_i|$  and  $e_i = |E_i|$ . In both generated graphs, we assume that  $n_i \approx \frac{n-r}{\sigma}$  and  $e_i \approx \frac{e}{\sigma}$  where  $r$  is the number of nodes without incoming edges.

We say a Wheeler graph  $G$  is *complete* if no more edges can be added to  $G$  while maintaining the Wheeler graph properties.



**Property 1.** Given number of nodes  $n$  and number of labels  $\sigma$ , the number of edges of a Wheeler graph is upper bounded by  $e_{max}$  where

$$e_{max} = n \times \sigma + n - \sigma - r = (n - 1)(\sigma + 1) - r + 1 \quad (5)$$

*Proof.* Consider the bipartite representation of a Wheeler graph  $G$  with number of nodes  $n$  and number of labels  $\sigma$ . Note that

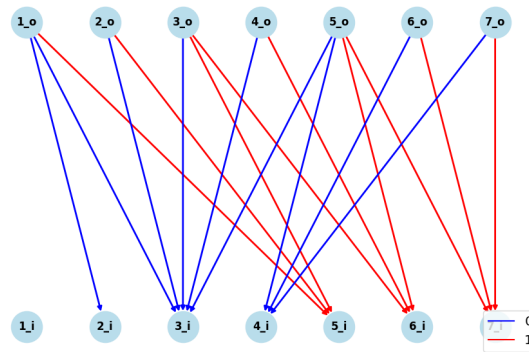
$$\sum_{i=1}^{\sigma} n_i = n - r. \quad (6)$$

Observe that for each label  $i$ , the number of edges that is labeled  $i$  is at most  $n + n_i - 1$ . Taking the sum of edges of each label and applying Equation (6), we have

$$e_{max} = \sum_{i=1}^{\sigma} (n + n_i - 1) = n \times \sigma - \sigma + \sum_{i=1}^{\sigma} n_i = n \times \sigma - \sigma + n - r. \quad (7)$$

□

One way of generating complete Wheeler graphs is to have all nodes connect to the first node of  $N_i$  and the last node additionally connect to the rest of the nodes in  $N_i$  for each label  $i$  (the last node has  $n_i$  outgoing edges in total). By randomly selecting  $n_i - 1$  nodes from  $N$  and connecting the selected nodes to consecutive nodes in  $N_i$ , a new complete Wheeler graph can be generated by appropriately shifting the destination node of each edge such that the Wheeler graph property is maintained. Figure 1 shows an example of a complete Wheeler graph with  $(n, e, \sigma, r) = (7, 18, 2, 1)$ . With a complete Wheeler graph of  $n$  nodes and  $\sigma$  labels, we are able to generate random Wheeler graphs with  $e < e_{max}$  edges by sampling  $e$  distinct edges from the complete Wheeler graph.



**Figure 1:** Complete Wheeler graph with  $(n, e, \sigma, r) = (7, 18, 2, 1)$ . In this example we have  $n_0 = n_1 = 3$ . The selected nodes for label 0 is node 1 and 5 and for label 1 node 3 and 5.

For generation of d-NFA Wheeler graphs, let  $x_k$  be the number of nodes with  $k$  outgoing edges of the same label  $i$ . Thus, given  $e$  and  $\sigma$  we have

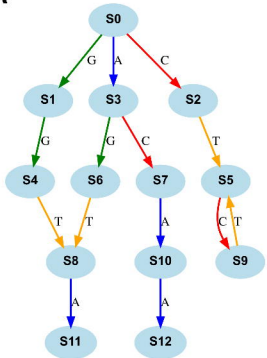
$$\sum_{k=1}^d k \cdot x_k = e_i \quad (8)$$

Also the minimum number of nodes  $n_{min}$  needed to accommodate  $e_i$  edges must be less than  $n_i$ . Note that  $x_1$  does not present in the formula below.

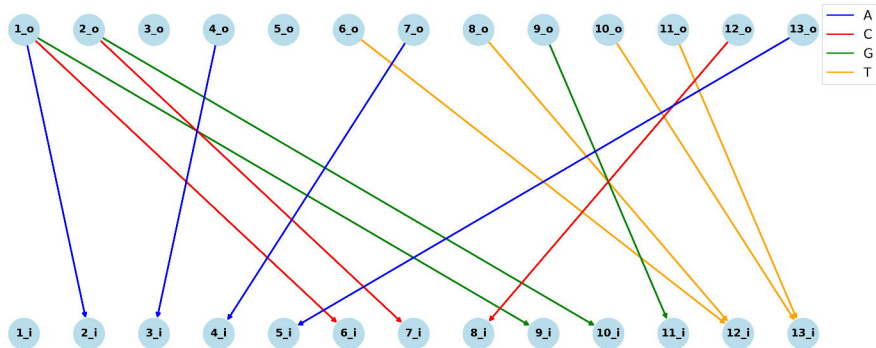
$$n_{min} = 1 + \sum_{k=2}^d (k-1) \cdot x_k \leq n_i \quad (9)$$

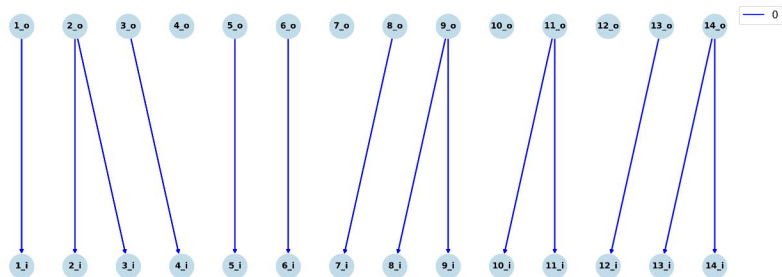
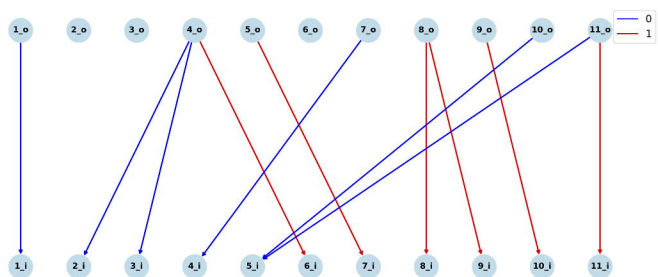
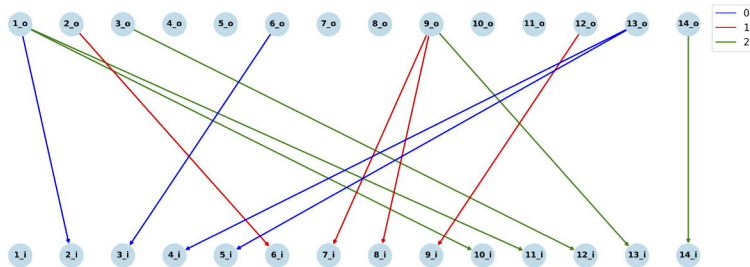
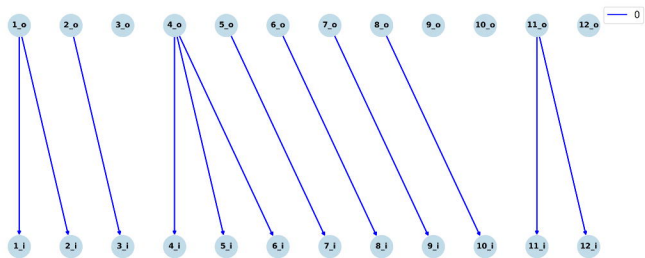
By finding a solution to  $x_k$  that satisfies Equation (8) and (9), we can guarantee a d-NFA Wheeler graph by randomly selecting  $x_k$  nodes and connecting  $k$  outgoing edges of label  $i$  to nodes in  $N_i$ . To generate non-trivial Wheeler graphs, we set all  $x_k$  to be the same, and if not possible assign the residual to  $x_1$  to satisfy Equation (8).



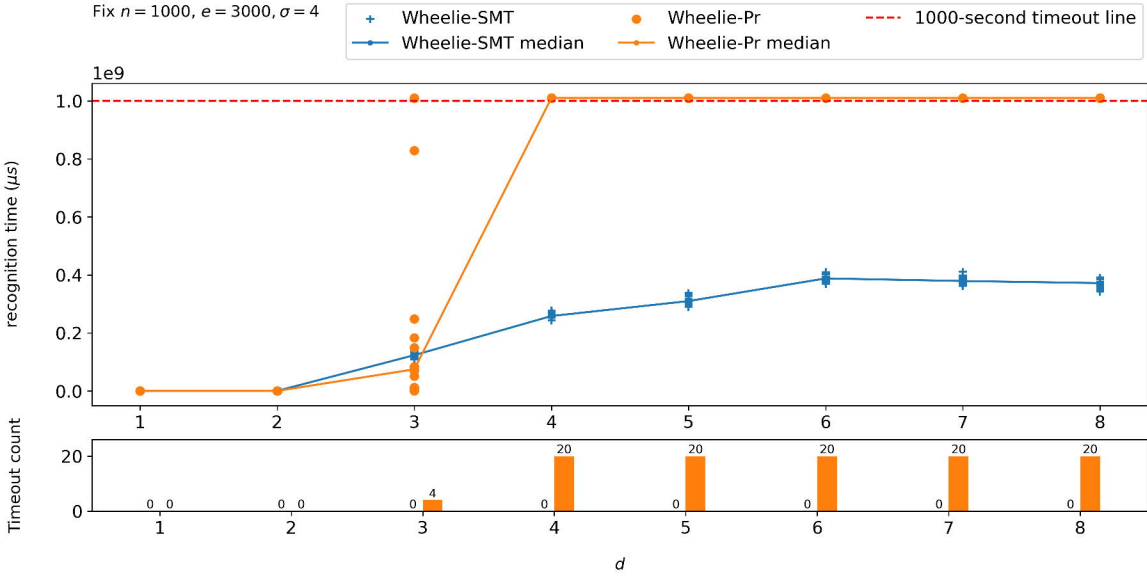
**A****B**

Node label	Wheeler order
S0	1
S3	2
S12	3
S10	4
S11	5
S2	6
S7	7
S9	8
S1	9
S6	10
S4	11
S5	12
S8	13

**C**

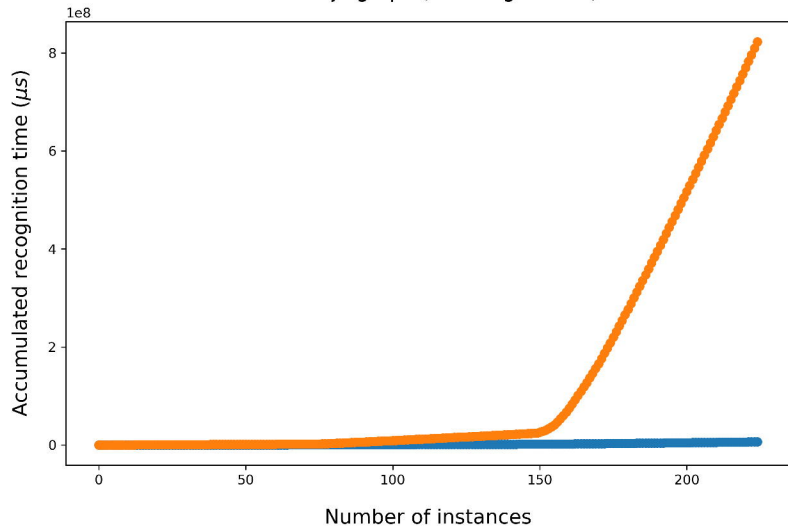
**A****B****C****D**

Fix  $n = 1000$ ,  $e = 3000$ ,  $\sigma = 4$







**A** De Bruijn graph (DNA alignments)**B** Reverse deterministic graph (DNA alignments)