

Canonicalization of Data Structures

Jakob Lykke Andersen

Research Group Bioinformatics and Computational Biology
Faculty of Computer Science
University of Vienna

Institute for Theoretical Chemistry
University of Vienna

Bled, February 2018

Introduction

Model

A mathematical object in some class M .

Example: a rational number, $\frac{3}{4}$

Introduction

Model

A mathematical object in some class M .

Example: a rational number, $\frac{3}{4}$

Representation

An object of an abstract data type R used to store the model.

Example: a pair of integers, $(3, 4)$

Introduction

Model

A mathematical object in some class M .

Example: a rational number, $\frac{3}{4}$

Representation

An object of an abstract data type R used to store the model.

Example: a pair of integers, $(3, 4)$

Implementation

An object of a concrete type used to store the model.

Example: `std::pair<int, int>(3, 4)`

Introduction

Model

A mathematical object in some class M .

Example: a rational number, $\frac{3}{4}$

Representation

An object of an abstract data type R used to store the model.

Example: a pair of integers, $(3, 4)$

Implementation

An object of a concrete type used to store the model.

Example: `std::pair<int, int>(3, 4)`

What if a model has multiple representations?

Example

$M = \mathbb{Q}$, the rational numbers

$R = \mathbb{Z} \times \mathbb{Z}$, pairs of integers

But $\frac{2}{5} = \frac{4}{10}$, so $(2, 5)$ should be considered “equal” to $(4, 10)$.

Notation: $(2, 5) \cong (4, 10)$ (“isomorphic to”)

$(2, 5) \stackrel{r}{=} (2, 5)$ (“representationally equal to”)

Canonicalization

Given a representation $G \in R$ find a new representation $C(G)$, such that:

- ▶ It represents the same model: $C(G) \cong G$
- ▶ All canonicalized isomorphic representations are the same:
 $\forall G' \in R, G' \cong G : C(G') \stackrel{r}{=} C(G)$

How do we specify and implement canonicalization in practice?

Representations

Besides the \cong^r operation we need:

- ▶ A class of operations, OP , that do not change the model.
- ▶ A total order $<^r$ among (isomorphic) representations.

Fraction Example:

OP :

- ▶ Multiplying with an integer: $(2, 5) \cdot 2 = (2 \cdot 2, 5 \cdot 5) \cong (2, 5)$
- ▶ Dividing with a common factor: $\frac{(4,10)}{2} = \left(\frac{4}{2}, \frac{10}{2}\right) \cong (4, 10)$
- ▶ (and compositions of those operations)

$<^r$:

- ▶ Prefer both positive over both negative: $(2, 5) <^r (-2, -5)$
- ▶ Prefer (neg., pos.) over (pos., neg.): $(-2, 5) <^r (2, -5)$
- ▶ Prefer smaller (absolute) numbers (lexicographically):
 $(2, 5) <^r (4, 10), \quad (1, 2) <^r (2, 3)$

Canonicalization

Given $G \in R$:

- ▶ Find $op \in \text{OP}$ that minimizes $op(G)$, wrt. $\overset{r}{<}$
- ▶ Return $op(G)$ as the canonical form.

Fraction Example:

Given (a, b) ,

- ▶ Find $f = \text{GCD}(|a|, |b|)$
- ▶ If $b < 0$: let $op = \text{DIV}(f) \circ \text{MUL}(-1)$
else: let $op = \text{DIV}(f)$
- ▶ Return $op((a, b))$

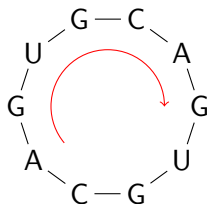
In Practice:

- ▶ Probably return op . The user can compute $op(G)$ if needed.
- ▶ $\overset{r}{<}$ may be implicitly defined by the canonicalization algorithm.

Example: Circular RNA (circRNA)

Representation: A sequence of symbols A, C, G, U.

Example: AGUGCAGUGC



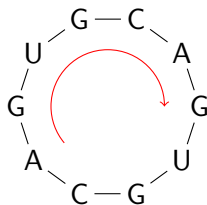
Operations: $\text{ROTATE}(i)$, for $i \in \mathbb{Z}$

Example: $\text{ROTATE}(2, \text{AGUGCAGUGC}) = \text{UGCAGUGCAG}$

$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Canonicalization: find the lexicographically smallest rotation
(can be done in linear time)

Example: Circular RNA (circRNA)



Representation: A sequence of symbols A, C, G, U.

Example: AGUGCAGUGC

Operations: $\text{ROTATE}(i)$, for $i \in \mathbb{Z}$

Example: $\text{ROTATE}(2, \text{AGUGCAGUGC}) = \text{UGCAGUGCAG}$

$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Canonicalization: find the lexicographically smallest rotation
(can be done in linear time)

Symmetry Discovery: op is a symmetry if $op(G) \stackrel{r}{=} G$

Example: $\text{ROTATE}(5)$ is a symmetry of AGUGCAGUGC, because
 $\text{ROTATE}(5, \text{AGUGCAGUGC}) = \text{AGUGCAGUGC}$
 $\stackrel{r}{=} \text{AGUGCAGUGC}$

$\text{ROTATE}(0)$ is a trivial symmetry

Example: Double Stranded RNA

Representation:

A pair of sequences of symbols A, C, G, U, of equal length.

Example: $\begin{matrix} \text{AGUGC} \\ \text{UCACG} \end{matrix}$

Operations: REVERSE \circ SWAP

Example: $(\text{REVERSE} \circ \text{SWAP}) \left(\begin{matrix} \text{AGUGC} \\ \text{UCACG} \end{matrix} \right) = \begin{matrix} \text{GCACU} \\ \text{CGUGA} \end{matrix}$

$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Example: $\begin{matrix} \text{AGUGC} \\ \text{UCACG} \end{matrix} \stackrel{r}{<} \begin{matrix} \text{GCACU} \\ \text{CGUGA} \end{matrix}$

Canonicalization: take the $\stackrel{r}{<}$ -smallest of the two possibilities

Example: Double Stranded RNA, Only Binding Structure

Swapping all A with U and G with C preserves structure.

Representation:

A pair of sequences of symbols A, C, G, U, of equal length.

Operations: REVERSE \circ SWAP and INVERT (= SWAP)

Example: $\text{INVERT} \begin{pmatrix} \text{AGUGC} \\ \text{UCACG} \end{pmatrix} = \begin{pmatrix} \text{UCACG} \\ \text{AGUGC} \end{pmatrix}$

$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Canonicalization: take the $\stackrel{r}{<}$ -smallest of the four possibilities

Example: Anti-Parallel Strong Traces, Take 1

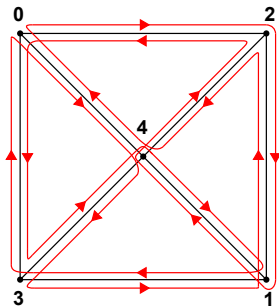
Model: A graph G (representing a polygon), with a closed walk visiting all edges twice and \langle more constraints \rangle .

Representation: A sequence of vertices $t = (v_{i_1}, v_{i_2}, \dots, v_{i_{2m}})$.

Operations: $\text{REVERSE}(t)$, $\text{ROTATE}(i, t)$, and $\text{PERMUTE}(\gamma, t)$ for any automorphism (i.e., symmetry) γ of G .

$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Canonicalization: take the $\stackrel{r}{<}$ -smallest (not trivial to do efficiently)



Example: Anti-Parallel Strong Traces, Take 1

Model: A graph G (representing a polygon), with a closed walk visiting all edges twice and \langle more constraints \rangle .

Representation: A sequence of vertices $t = (v_{i_1}, v_{i_2}, \dots, v_{i_m})$.

Operations: REVERSE(t), ROTATE(i, t), and PERMUTE(γ, t) for any automorphism (i.e., symmetry) γ of G .

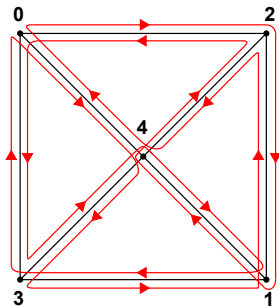
$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Canonicalization: take the $\stackrel{r}{<}$ -smallest (not trivial to do efficiently)

But what about the graph?

What is a vertex?

What is the representation?



Example: Graphs, Part 1

Model: A graph $G = (V, E)$.

Representation: An adjacency list which implicitly assigns $1, 2, \dots, n$ to V .

Operations: $\text{PERMUTE}(\gamma)$ for any permutation of $1, 2, \dots, n$.

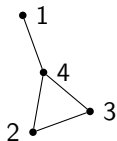
$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Canonicalization: ⟨more on this later⟩

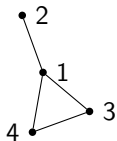
Graph Representation and Graph Permutation

$$G = (V, E) \quad V = \{1, 2, \dots, n\}$$

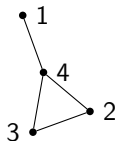
Isomorphic graphs, different representations:



G_1



G



G_2

Adjacency list representation (with sorted neighbour lists):

1: 4
2: 3, 4
3: 2, 4
4: 1, 2, 3

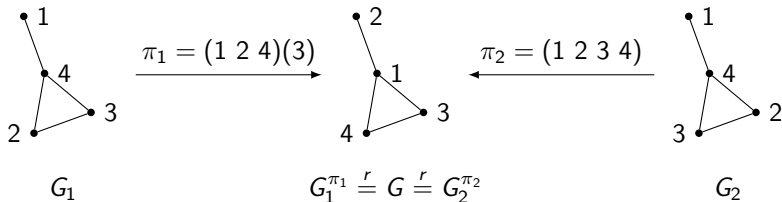
1: 2, 3, 4
2: 1
3: 1, 4
4: 1, 3

1: 4
2: 3, 4
3: 2, 4
4: 1, 2, 3

Graph Representation and Graph Permutation

$$G = (V, E) \quad V = \{1, 2, \dots, n\}$$

Isomorphic graphs, different representations:



Adjacency list representation (with sorted neighbour lists):

1: 4
2: 3, 4
3: 2, 4
4: 1, 2, 3

1: 2, 3, 4
2: 1
3: 1, 4
4: 1, 3

1: 4
2: 3, 4
3: 2, 4
4: 1, 2, 3

Example: Anti-Parallel Strong Traces, Take 2

Model: A graph G (representing a polygon), with a closed walk visiting all edges twice and \langle more constraints \rangle .

Representation: An adjacency list, and a sequence of integers $t = (v_{i_1}, v_{i_2}, \dots, v_{i_{2m}})$.

Operations:

- ▶ $\text{REVERSE}(t)$
- ▶ $\text{ROTATE}(i, t)$
- ▶ $\text{PERMUTE}(\gamma, t)$ for any automorphism γ of G .
- ▶ $\text{PERMUTE}(\gamma, t, G)$ for any permutation γ of V .

$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Canonicalization: take the $\stackrel{r}{<}$ -smallest

Example: Anti-Parallel Strong Traces, Take 3

Model: A graph G (representing a polygon), with a closed walk visiting all edges twice and \langle more constraints \rangle .

Representation: An adjacency list, and a sequence of integers representing a *gap vector* $g = (a_1, a_2, \dots, a_{2m})$.

Operations:

- ▶ $\text{ROTATE}(i, g)$
- ▶ $\text{MAKEGAP} \circ \text{REVERSE} \circ \text{MAKETRACE}$

$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Canonicalization: find the $\stackrel{r}{<}$ -smallest rotation of g and its reverse.

Example: Graphs, Continued

Model: A graph $G = (V, E)$.

Representation: An adjacency list which implicitly assigns $1, 2, \dots, n$ to V .

Operations: $\text{PERMUTE}(\gamma)$ for any permutation of $1, 2, \dots, n$.

$\stackrel{r}{=}$ and $\stackrel{r}{<}$: component-wise and lexicographic comparison

Computational Complexity: $\exp(O(\sqrt{n \log n}))$

Brute-Force Algorithm:

1. Construct G^γ for all permutations $\gamma \in S_n$.
2. Select the “best” one (for example the $\stackrel{r}{<}$ -smallest).

[Babai and Luks, STOC, 1983]

[Babai, Handbook of Combinatorics, 1996]

Existing Tools for Canonicalization in Practice

Published Tools: nauty, Traces, Bliss (and Saucy and Conauto)

- ▶ All based on the idea of [individualization-refinement](#).
- ▶ Different sets of heuristics and variations.
- ▶ Many more algorithm variations are possible.
- ▶ Which is the best? for a specific class of graphs?
- ▶ What if the graph has vertex and edge labels?
- ▶ What if those labels are “complicated”? (e.g., stereo-info)

[McKay, Congressus Numerantium, 1981] [McKay and Piperno, J. Symb. Comp., 2014] [Junttila and Kaski, ALENEX, 2007] [Darga et al., DAC, 2008] [López-Presa and Fernández Anta, SEA, 2009]

Existing Tools for Canonicalization in Practice

Published Tools: nauty, Traces, Bliss (and Saucy and Conauto)

- ▶ All based on the idea of **individualization-refinement**.
- ▶ Different sets of heuristics and variations.
- ▶ Many more algorithm variations are possible.
- ▶ Which is the best? for a specific class of graphs?
- ▶ What if the graph has vertex and edge labels?
- ▶ What if those labels are “complicated”? (e.g., stereo-info)

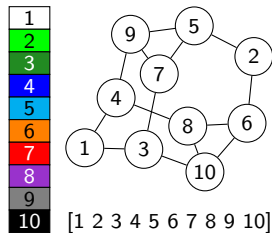
GraphCanon: [Andersen and Merkle, ALENEX, 2018]

- ▶ A generic C++ library for canonization algorithms.
- ▶ Algorithm variations implementable as individual plugins.
- ▶ Allows direct comparison of algorithm variations.
- ▶ Lower barrier of entry for implementing new ideas.
- ▶ Generality wrt. vertex/edge attributes.

[McKay, Congressus Numerantium, 1981] [McKay and Piperno, J. Symb. Comp., 2014] [Junttila and Kaski, ALENEX, 2007] [Darga et al., DAC, 2008] [López-Presa and Fernández Anta, SEA, 2009]

Individualization-Refinement Paradigm

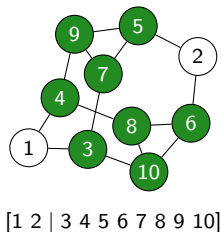
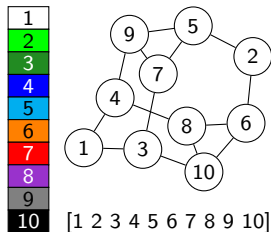
Initially: all vertices are unordered (same colour).



Individualization-Refinement Paradigm

Refine the ordering by propagation of “cheap” local information.

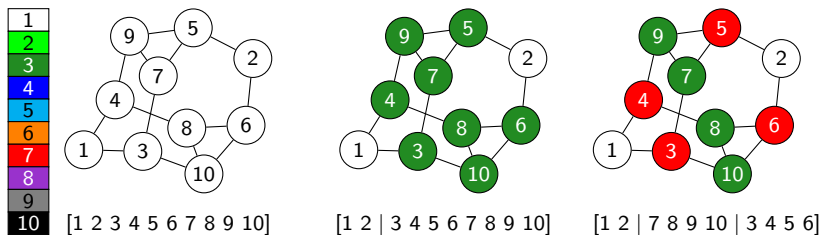
Example: sort and partition by degree (1D Weisfeiler-Leman).



Individualization-Refinement Paradigm

Refine the ordering by propagation of “cheap” local information.

Example: sort and partition by degree (1D Weisfeiler-Leman).



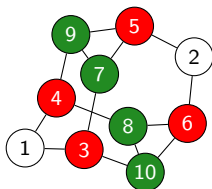
Individualization-Refinement Paradigm

Let this be the root of a search tree, and select a colour.

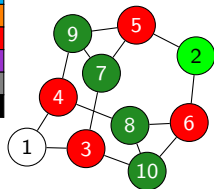
For each vertex of that colour;

create a child with this vertex given a unique new colour.

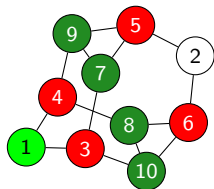
1
2
3
4
5
6
7
8
9
10



[1 2 | 7 8 9 10 | 3 4 5 6]

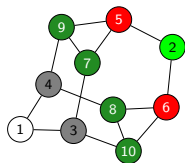


[1 | 2 | 7 8 9 10 | 3 4 5 6]

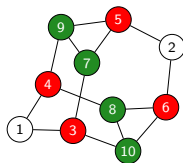


[2 | 1 | 7 8 9 10 | 3 4 5 6]

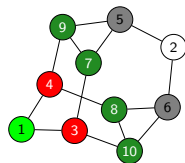
Individualization-Refinement Paradigm



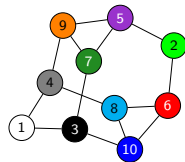
$$\pi_{(1)} = [1 \mid 2 \mid 7 \ 8 \ 9 \ 10 \mid 5 \ 6 \mid 3 \ 4]$$



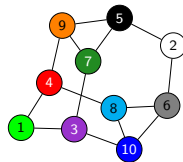
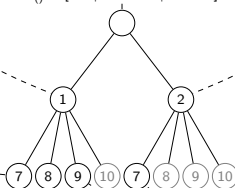
$$\pi_0 = [1 \ 2 \mid 7 \ 8 \ 9 \ 10 \mid 3 \ 4 \ 5 \ 6]$$



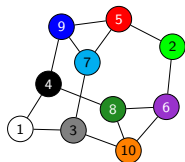
$$\pi_{(2)} = [2 \mid 1 \mid 7 \ 8 \ 9 \ 10 \mid 3 \ 4 \mid 5 \ 6]$$



$$\pi_{(1,7)} = [1 \mid 2 \mid 7 \mid 10 \mid 8 \mid 9 \mid 6 \mid 5 \mid 4 \mid 3]$$



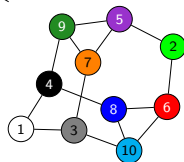
$$\pi_{(2,7)} = [2 \mid 1 \mid 7 \mid 10 \mid 8 \mid 9 \mid 4 \mid 3 \mid 6 \mid 5]$$



$$\pi_{(1,8)} = [1 \mid 2 \mid 8 \mid 9 \mid 7 \mid 10 \mid 5 \mid 6 \mid 3 \mid 4]$$

1
2
3
4
5
6
7
8
9
10

Colour order



$$\pi_{(1,9)} = [1 \mid 2 \mid 9 \mid 8 \mid 10 \mid 7 \mid 6 \mid 5 \mid 3 \mid 4]$$

Algorithm Variation

Categories

- ▶ Tree traversal
- ▶ Target cell selection
- ▶ Refinement
- ▶ Pruning with automorphisms
- ▶ Detection of implicit automorphisms
- ▶ Node invariants

GraphCanon: A common extension infrastructure.

Each variation implemented as a **visitor**:

- ▶ A set of callback methods for events of interest.
- ▶ Additional data structures instantiated
 - ▶ per search tree
 - ▶ per tree node

Benchmarks

44 graph collections, with 4,715 graphs in total.

Time limit: 1000 s

Memory limit: 8 GB

Repetitions: 5

Algorithm configurations: $\{\text{BFSExp}, \text{DFS}\} \times \{\text{F}, \text{FL}, \text{FLM}\}$

Compute nodes with two Intel E5-2680v3 CPUs (24 cores)

Compute node hours: approx. 12,000

BFSExp with FLM is often best.

CFI-Rigid: [Neuen and Schweitzer, ESA, 2017]

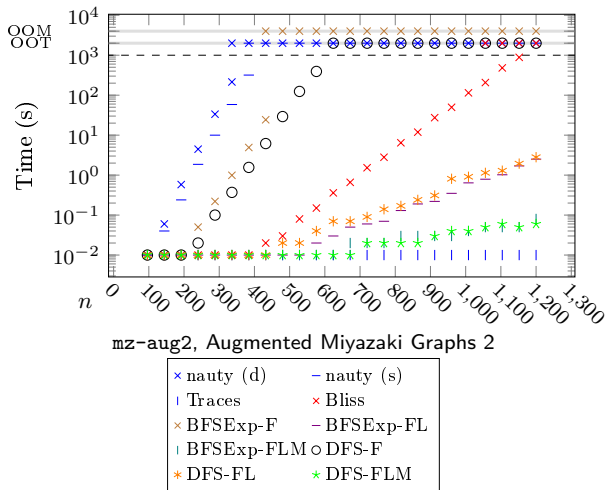
nauty, Traces: [<http://pallini.di.uniroma1.it/Graphs.html>]

Bliss: [<http://www.tcs.hut.fi/Software/bliss/benchmarks/index.shtml>]

Conauto: [<https://sites.google.com/site/giconauto/home/benchmarks>]

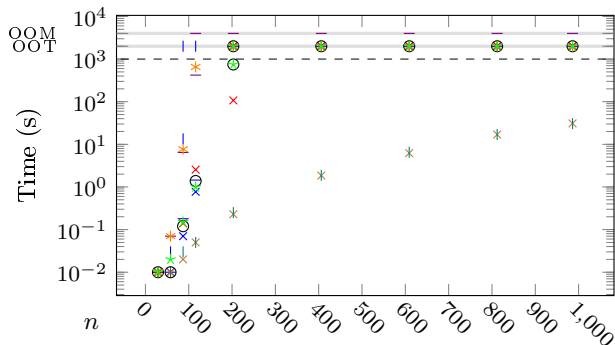
Saucy: [<http://vlsicad.eecs.umich.edu/BK/SAUCY/>]

Tree Traversal and Target Cell Selector

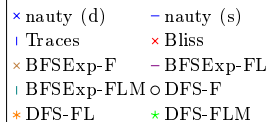


Similar characteristics observed for other Miyazaki graphs.

Tree Traversal and Target Cell Selector



usr, Union of Strongly Regular Graphs



CFI-Rigid

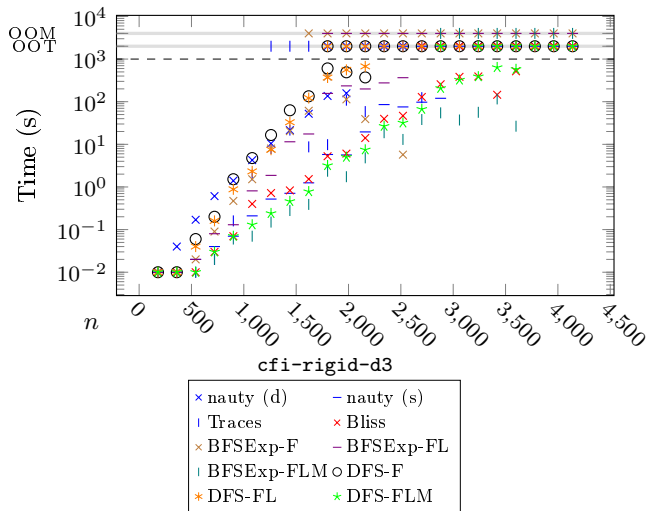
- ▶ 6 collections
- ▶ Designed to be the hard benchmarks.
- ▶ Expected to have very little symmetry.

Algorithm configurations:

$$\{\text{BFSExp, DFS}\} \times \{\text{F, FL, FLM}\} \times 2^{\{\text{PL, Q, T}\}}$$

Col.	Group	Reduction	Best Algorithm	Invariants Matter	FLM Sep.	Max. Solved	n
d3	D_3	—	BFSExp-FLM	yes (any)	yes	3,600	
z3	\mathbb{Z}_3	—	BFSExp-FLM	yes (any)	yes	3,780	
z2	\mathbb{Z}_2	—	Bliss, nauty (s)	yes (any)	no	2,992	
r2	\mathbb{Z}_2	R^*	Bliss, nauty (s)	no	no	1,584	
s2	\mathbb{Z}_2	B^*	FLM, Bliss, nauty (s)	yes (PL or Q)	no	2,496	
t2	\mathbb{Z}_2	$R^* \circ B^*$	FLM, Bliss, nauty (s)	yes (PL or Q)	yes	1,056	

CFI-Rigid



Summary

- ▶ Canonicalization is a general principle.
- ▶ The concepts can be applied to any data structure.
- ▶ Brute-force: make it a graph.

GraphCanon:

- ▶ Generic algorithm framework.
- ▶ (Relatively) easy to develop new variations.
- ▶ Allows direct comparison of algorithmic ideas.
- ▶ Competitive with established tools.
- ▶ https://github.com/jakobandersen/graph_canon
- ▶ Very easy to extract data for visualization:
https://jakobandersen.github.io/graph_canon_vis/

MØD v0.7 (to be released soonTM):

- ▶ Integrates GraphCanon .
- ▶ Finally, true canonical SMILES strings!
- ▶ The automorphism group of molecules is now available.
(important for atom tracing)

Algorithm Variation

Tree Traversal:

- ▶ `nauty`, `Bliss`: depth-first (DFS)
- ▶ `Traces`: breadth-first with experimental paths (BFSExp)
- ▶ `GraphCanon`:
 - ▶ Arbitrary traversals are possible.
 - ▶ Garbage collected search tree via reference counting.
 - ▶ Extensions must keep owning references to tree nodes.
 - ▶ Implemented: DFS, BFSExp, and a new hybrid (BFSExpM).

Target Cell Selector:

- ▶ Many have been developed.
- ▶ Currently implemented:
 - ▶ first (F)
 - ▶ first largest (FL)
 - ▶ first largest with maximum number of non-uniformly joined neighbour cells (FLM)

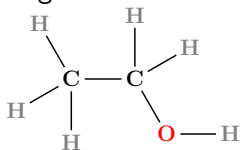
Algorithm Variation

Node Invariants:

- ▶ Totally ordered isomorphism-invariant information.
- ▶ Invariants can be implemented independently.
- ▶ A special visitor coordinates invariants.
- ▶ Implemented:
 - ▶ cell splitting positions (T), from Traces
 - ▶ quotient graph values (Q), from nauty, Traces, Bliss
(but not hashed)
 - ▶ partial leaf (PL), from Bliss (but not hashed)
Construct parts of the permuted graph earlier in the tree.

Refinement functions implemented:

- ▶ 1D Weisfeiler-Leman, generalized to exploit edge attributes.
- ▶ A function to handle degree-1 vertices.



Algorithm Variation

Detection of implicit automorphisms:

- ▶ Sometimes we can detect/guess automorphisms at internal tree nodes.
- ▶ **nauty**: several special cases of ordered partitions.
- ▶ **Saucy**: heuristics for guessing sparse automorphisms.
- ▶ **Traces**: reportedly a generalization of the Saucy heuristics.
- ▶ **Implemented**:
 - ▶ Partitions where all cells have size 1 or 2.
 - ▶ The degree-1 vertex refinement function.

Pruning with automorphisms:

Calculation of orbits in stabilizers of the found automorphisms.

Stabilizer calculation:

- ▶ **nauty** (early versions) and **Bliss**: conservative (implemented)
- ▶ **Traces** and **nauty** (recent versions): randomized Schreier-Sims

The implemented visitor for automorphism pruning is generic with respect to stabilizer implementation.