

TADASHI: Enabling ML to explore transformations with guaranteed correctness

E. Vatai, A. Drozd, I. R. Ivanov, J. E. Batista, Y. Ren, M. Wahib

Riken R-CCS, High Performance Artificial Intelligence Systems Research Team,
Japan

Outline

Intro

Motivation

`pip install tadashi`

The Long Term Vision

... and how it's goin'

Bonus 3 Slide Polyhedral Tutorial

The Team



Emil
VATAI



Aleksandr
DROZD



Ivan R.
IVANOV



João E.
BATISTA



Mohamed
WAHIB

THIS IS YOUR MACHINE LEARNING SYSTEM?

|
YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

|
WHAT IF THE ANSWERS ARE WRONG?
|

JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.



ML approaches to correctness in code generation

- ▶ Extensive unit testing
 - ▶ Coverage
 - ▶ Writing tests is not trivial
- ▶ Surrogate ML model
- ▶ Round-trip
 - ▶ Trust issues
- ▶ Limit ML to short sequences of instructions
- ▶ Limit ML to a set of determined transformations
 - ▶ Not general
- ▶ Symbolic execution
 - ▶ Not well established

Edsger W. Dijkstra

"Program testing can be used to show the presence of bugs, but never to show their absence!"

Some people don't even do that

Job interviewer: It says in your CV that you are quick at mathematics.
What is 17×19 ?

Me: 36

Interviewer: That's not even close
Me: But it was quick



TechCrunch
<https://techcrunch.com>

Microsoft CEO says up to 30% of the company's code was written by AI

29 Apr 2025 — Microsoft CEO Satya Nadella said that 20% to 30% of code inside the company's repositories was "written by software" — meaning...

Hacker News
<https://news.ycombinator.com>

Microsoft admits almost all major Windows 11 core features are broken

1 day ago — Except it's not totally terrible. What makes it terrible is the privacy invasions and the consent issues Microsoft has with its users, along ...



All you need is to sample the set of correct transformations

Scope:

- ▶ Hotspots

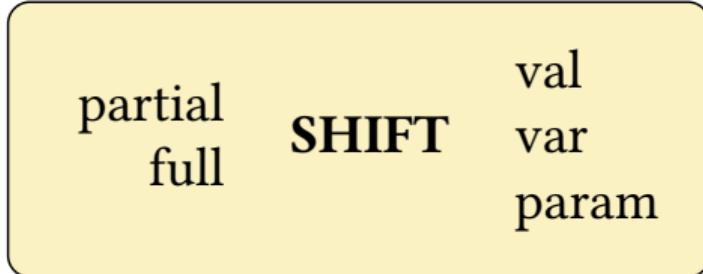
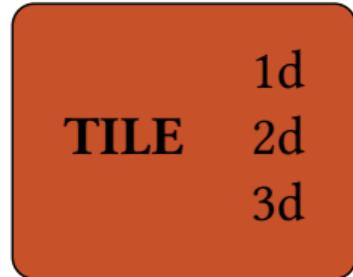
How it was done before

- ▶ By hand!
- ▶ Pluto, PPCG: Heuristics that don't scale
- ▶ Different approaches AI/Compiler/Applications people

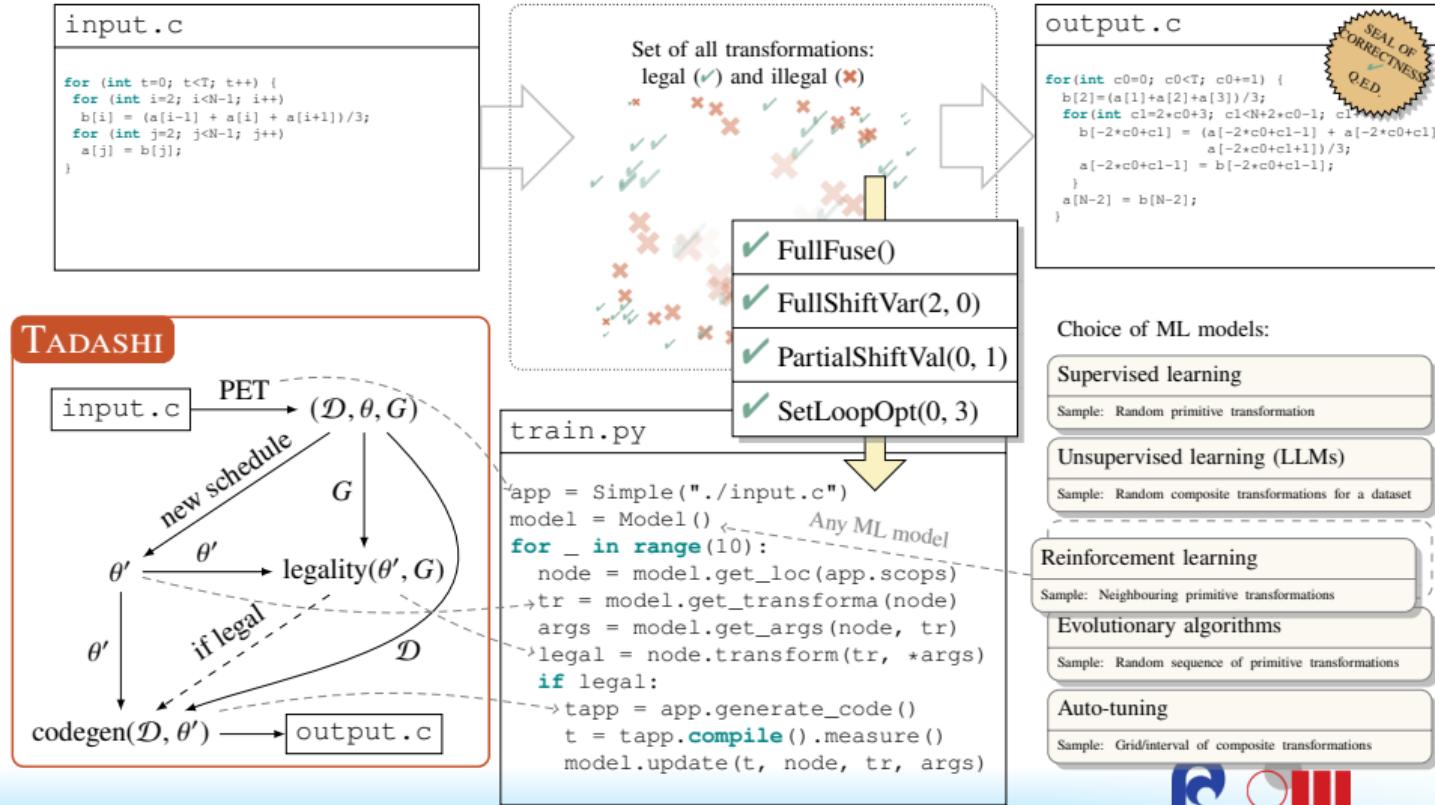
Our view: Not writing, but rewriting

- ▶ ML codegen is **transformations on a reference implementation**
- ▶ ML can **explore the space of transformations** to find a faster version
- ▶ But we also need **correctness**

Loop Transformations



TADASHI: loop transformations with correctness check

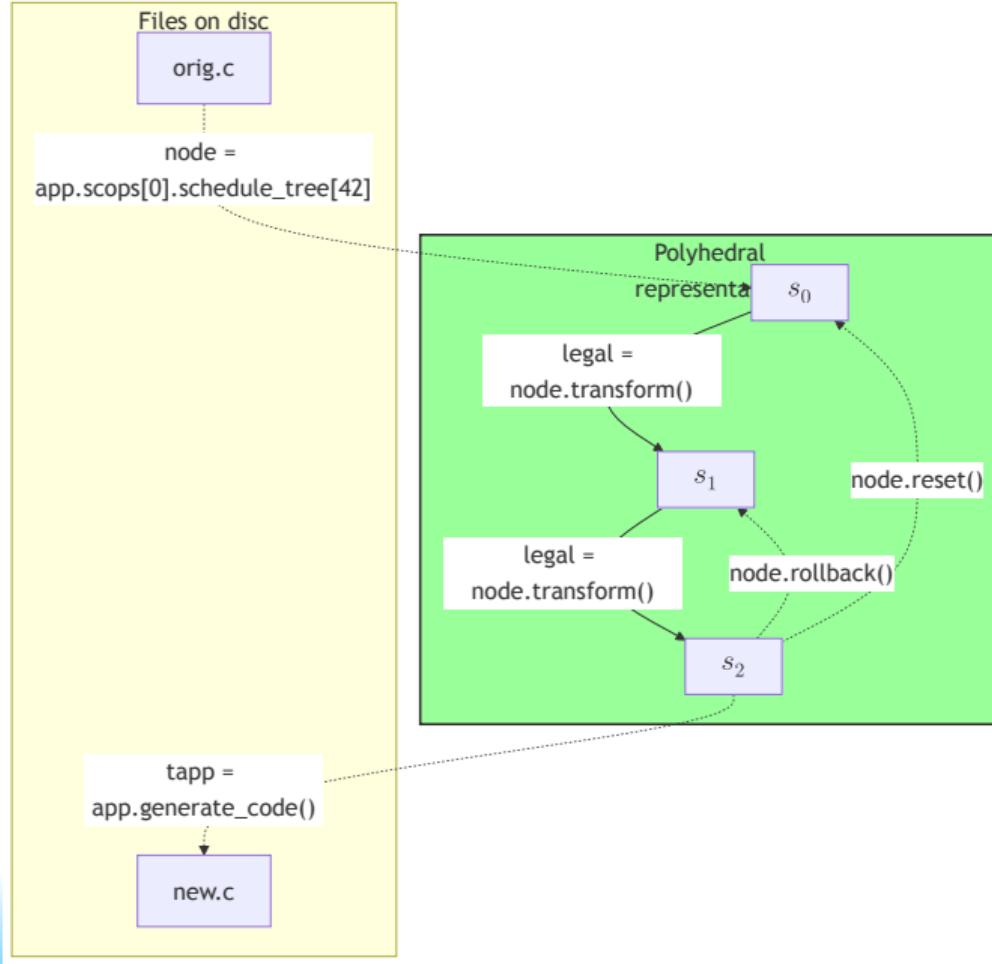


End-to-end example (bootstrapping for ML dev)

```
1 from pathlib import Path
2 from random import choice, seed
3 import tadashi
4 from tadashi.apps import Simple
5 seed(1234)
6 dir_path = Path(__file__).parent
7 examples_path = dir_path if dir_path.name == "examples" else "exa
8 app = Simple(f"{examples_path}/inputs/depnodedep.c")
9 node = app.scops[0].schedule_tree[1]
10 tr = choice(node.available_transformations)
11 args = choice(node.get_args(tr, -10, 10))
12 legal = node.transform(tr, *args)
13 app.compile()
14 transformed_app = app.generate_code()
15 transformed_app.compile()
16 print(f"{app.measure()=}")
17 print(f"{transformed_app.measure()=}")
```

Transformation list

```
1 trs = [[3, "FULL_SPLIT"],  
2         [1, "TILE2D", 32, 35],  
3         [5, "INTERCHANGE"]]
```

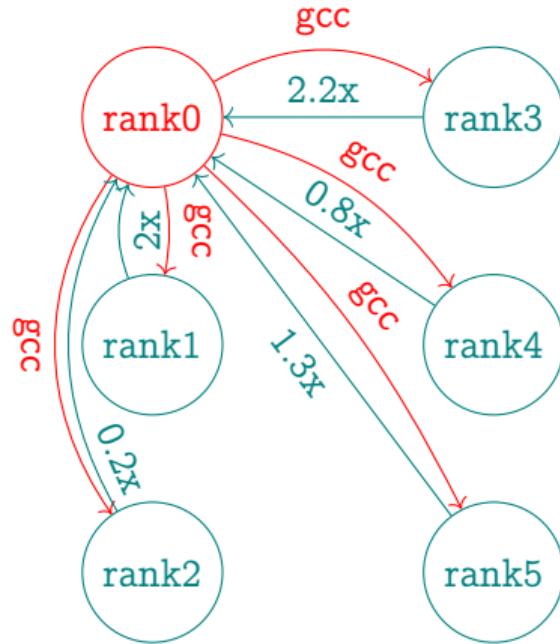


Benchmarking harness

The **compiler** spends a few minutes optimising the code on a single node.

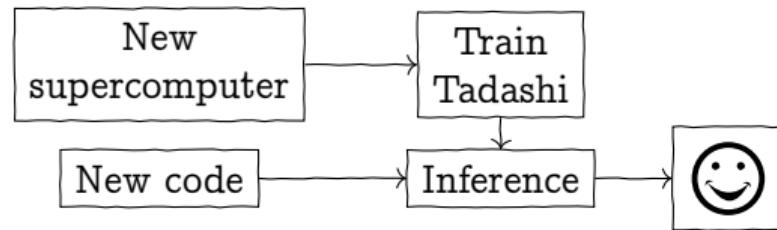
Why not use the entire supercomputer?

- ▶ Benchmarking harness
 - ▶ Distribute compiling and running of benchmarks
 - ▶ Collect the speedups to remote nodes
 - ▶ Uses **MPI4py**, and exposed to the user as **concurrent.futures** (from the Python standard)



Grand vision

Large scale optimisation framework

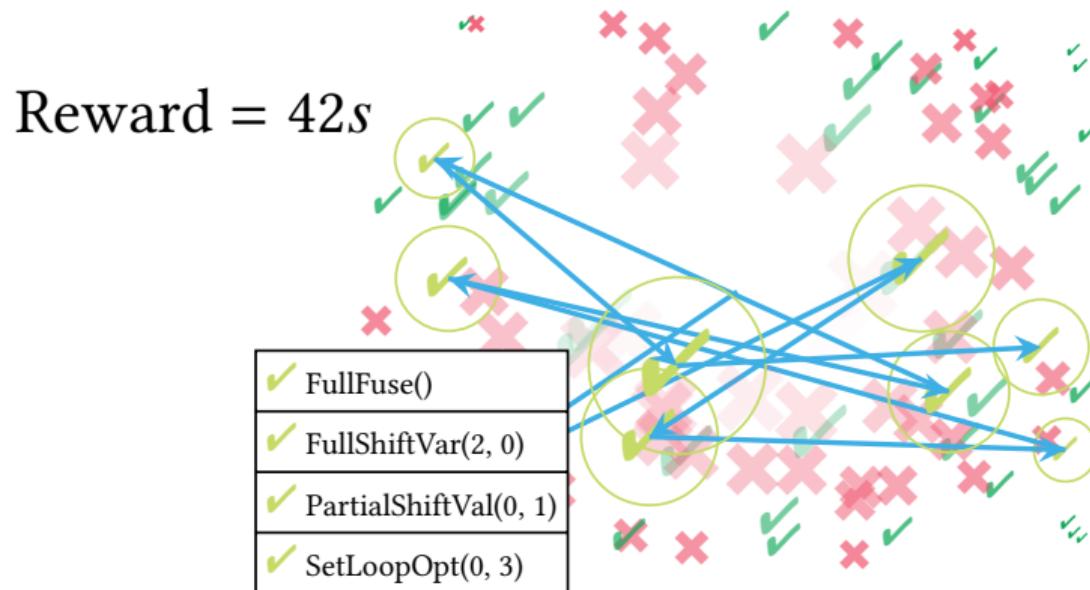


Wise words

Spending 6h on the whole machine to make the code 5% faster will pay off. – N. D.

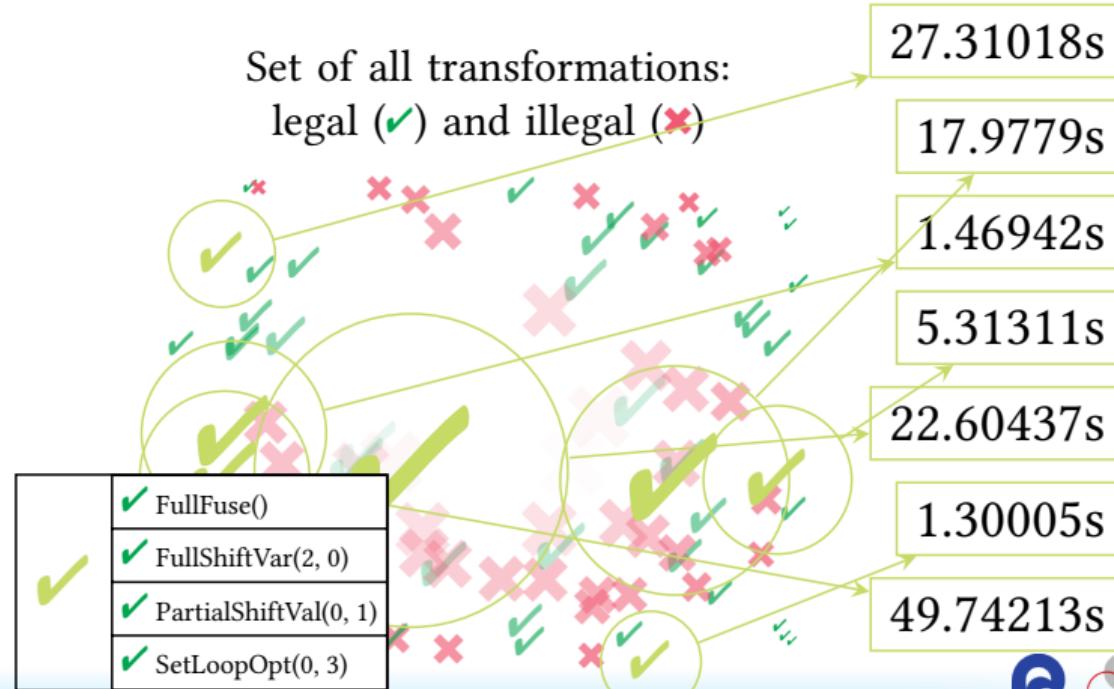
Reinforcement learning

- Actions = primitive transformations; Reward = walltime; Environment = correctness



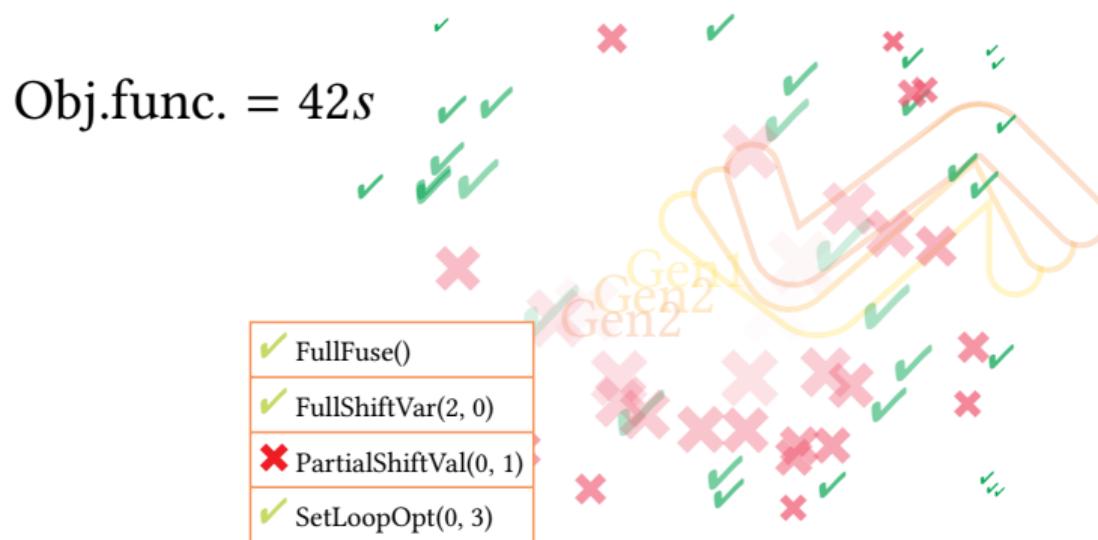
Supervised learning

- Dataset generation (legal transformations); Sample label: runtime



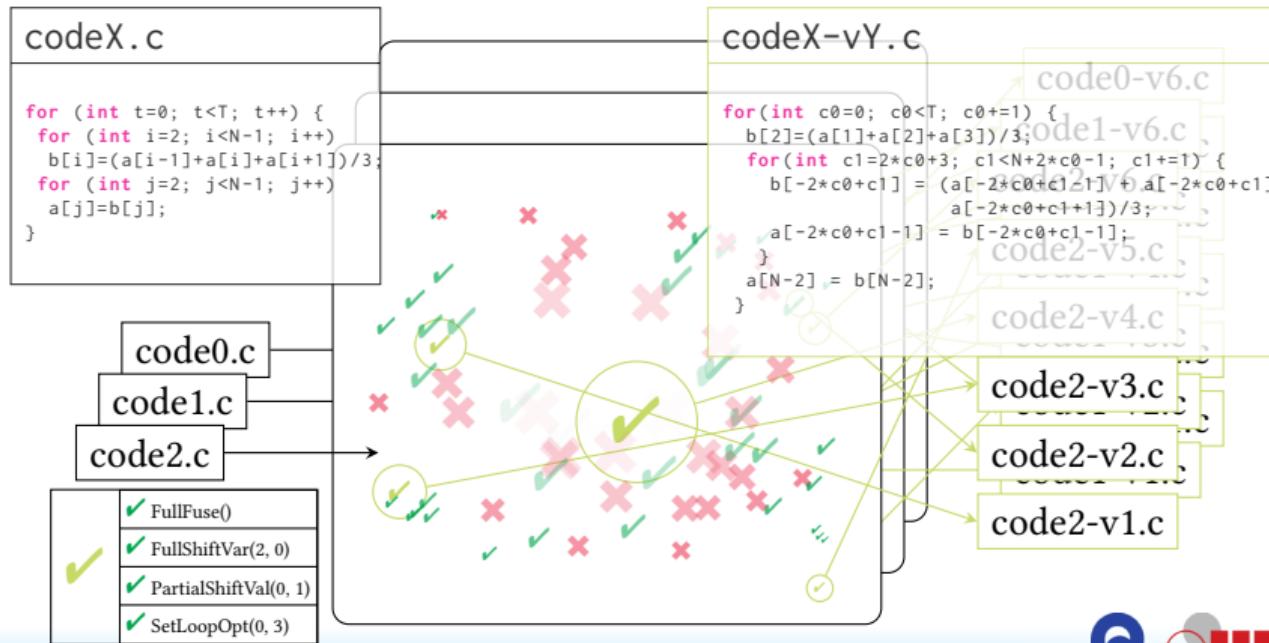
Evolutionary algorithms

- ▶ Exploration and explorations; Candidates = transformations; Objective function = runtime



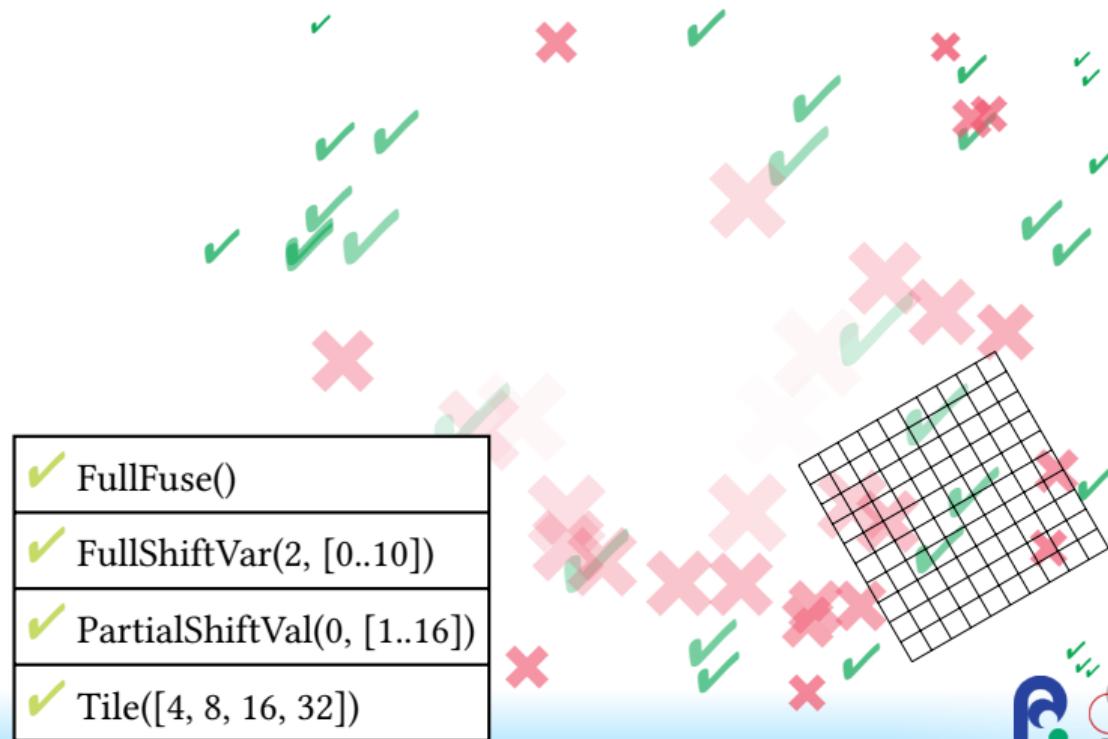
LLM agents

- ▶ Dataset generation; High quality correct transformations; Caveat!
Hallucinations are possible!



Auto-Tuning

- Brute forcing on a bounded region of the space (grid search)



ML opportunities

Different levels of abstraction/representations:

- ▶ source code [Stein: Paraphrasing etc.]
- ▶ abstract syntax tree (AST) [Shido: Tree-LSTM etc.]
- ▶ polyhedral [Baghdadi: Tiramisu]
- ▶ dependency graphs [Cummins: PrograML]
- ▶ intermediate Representations (IR) [Ben-nun: inst2vec]
- ▶ assembly instructions [Deepmind: Faster sorting]

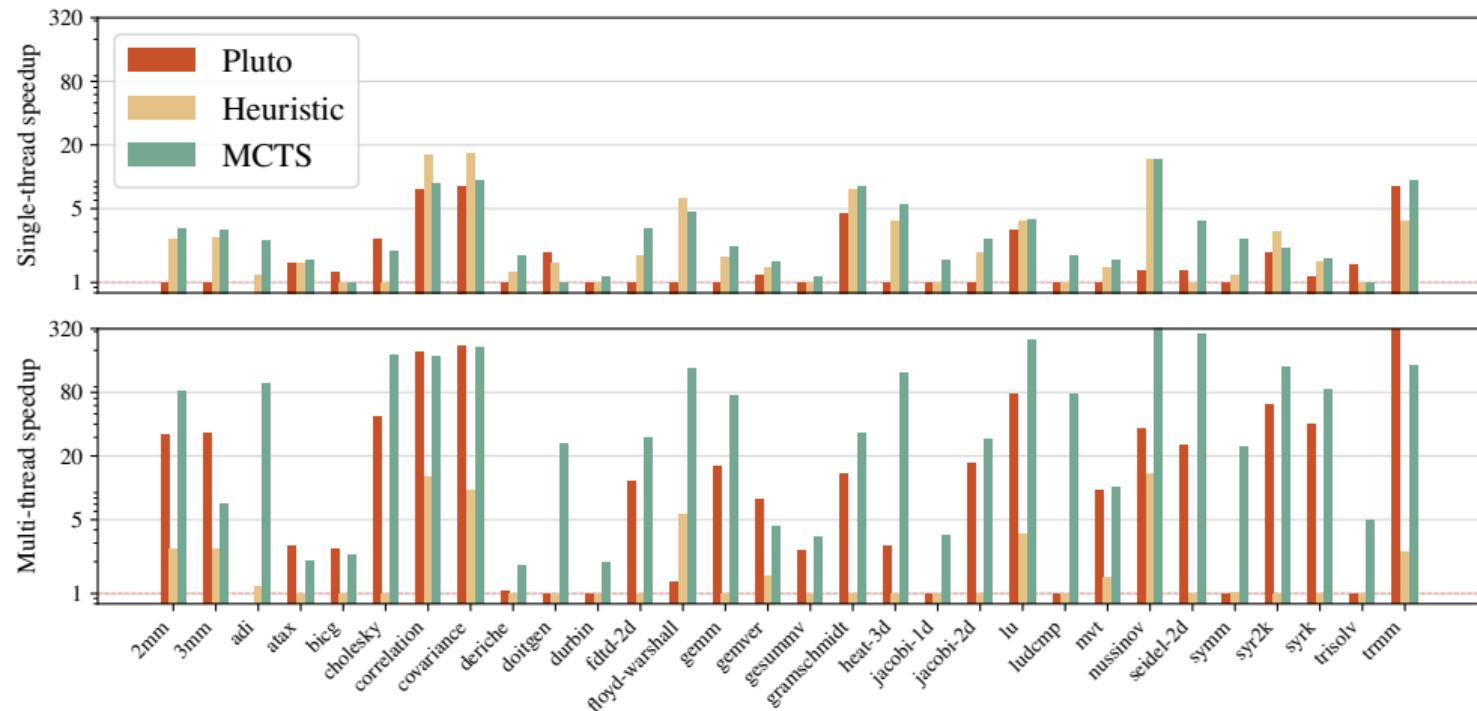
Three levels of transfer learning

1. **No Transfer**: Train from scratch for each (SW, HW) pair
2. Transfer to **new software**
3. Transfer to **new hardware**

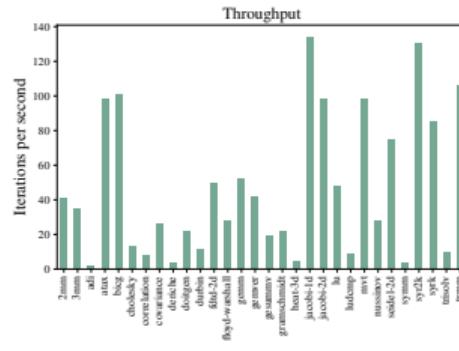
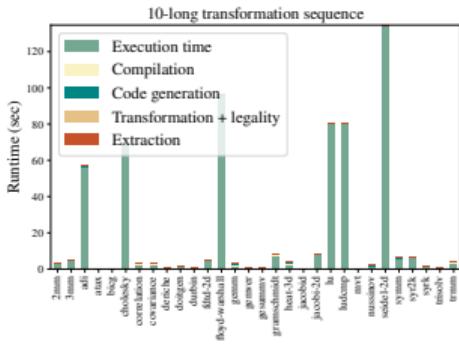
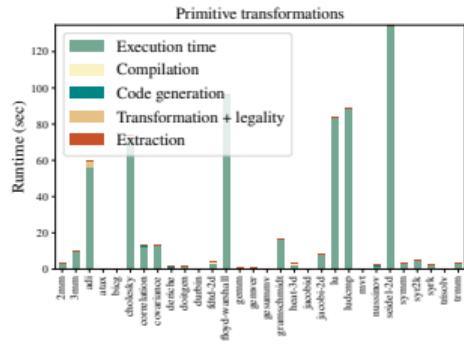
Less retraining, Better exploration \Rightarrow better results for a given kernel

Speedups: Heuristic, MCTS, GP

for heuristic, Monte Carlo Tree Search, Pluto



Overhead breakdown (primitive, 10 seq), Throughput



Real world applications

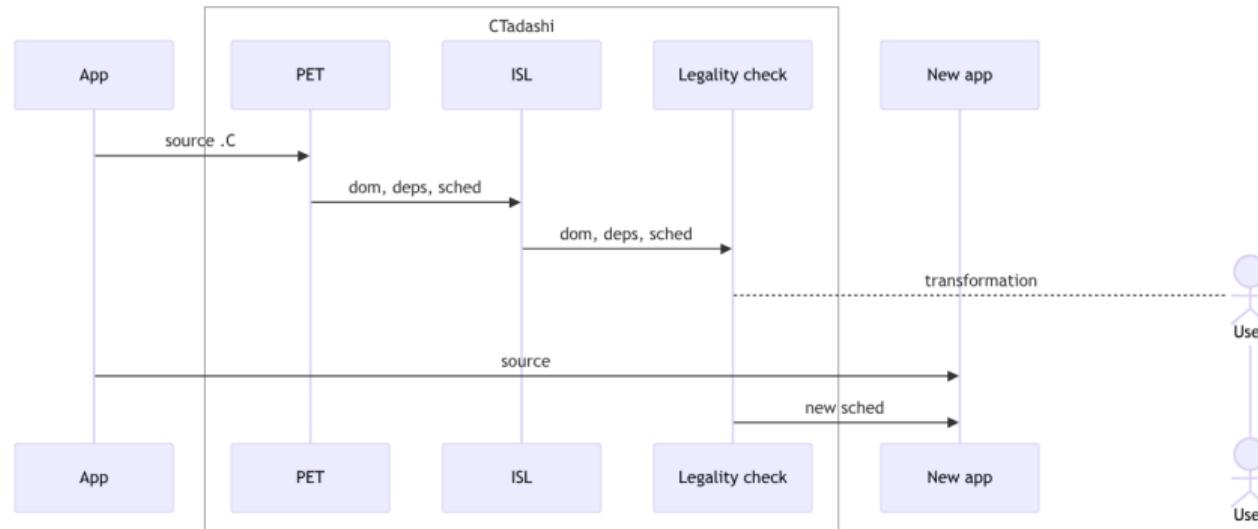
- ▶ miniAMR 5%
- ▶ Finite-difference time-domain FDTD 13x

Known limitations

- ▶ C only
- ▶ SCoP (polyhedral transformations)
- ▶ PET/ISL
- ▶ CPU only

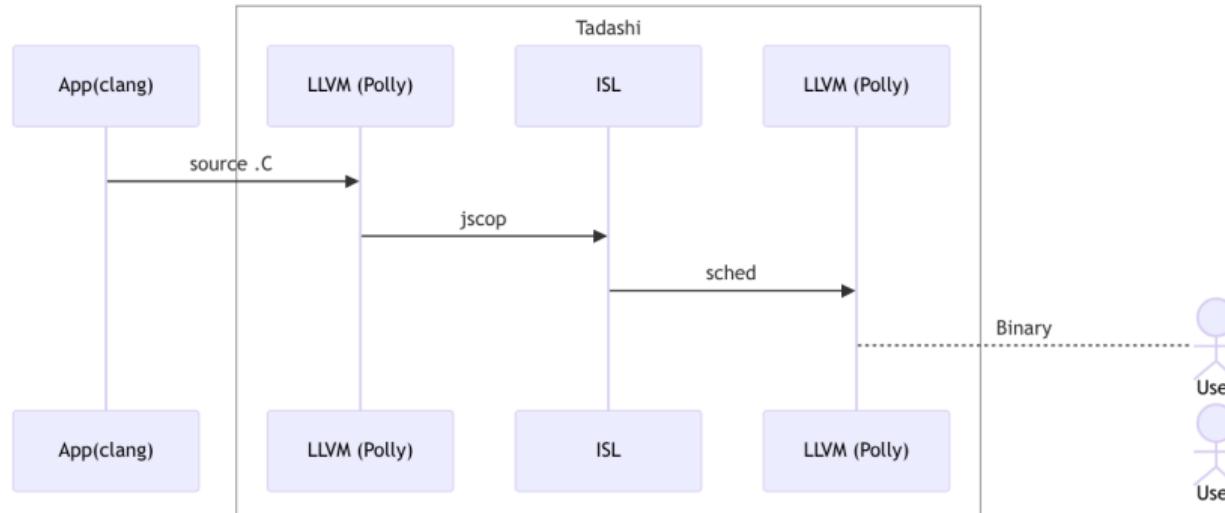
Future Works: Fortran

Current backend: PET+ISL



Future Works: Fortran

WIP backend: LLVM+ISL



Roadmap

- ▶ Codegen: GPU
- ▶ ML
- ▶ Beyond polyhedral

That's all folks!

- ▶ i. <https://vatai.github.io/>
- ▶ ii. <https://github.com/vatai/tadashi/>
- ▶ iii. <https://arxiv.org/abs/2410.03210>



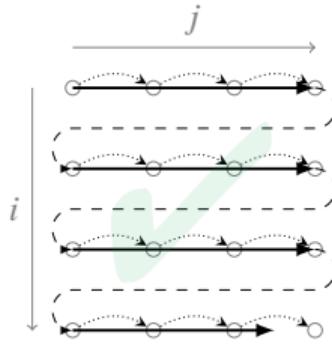
Polyhedral model

Components

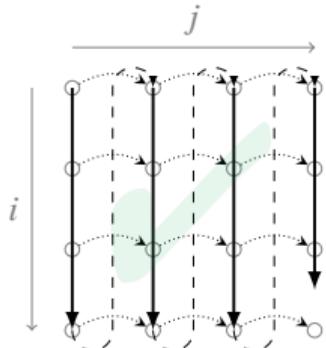
1. $\mathcal{D}_S = \{S[\vec{i}] \in \mathbb{Z}^n : \mathbf{A}\vec{i} + \mathbf{b} \leq 0\}$
2. $\theta(S[i, j]) = t = (i, j)$
3. $G = (V, E)$:
 - ▶ $V = \{S_0, S_1, \dots\}$,
 - ▶ $E = \{S_i[\vec{d}] \mapsto S_j[\vec{r}], \dots\}$

Mini example

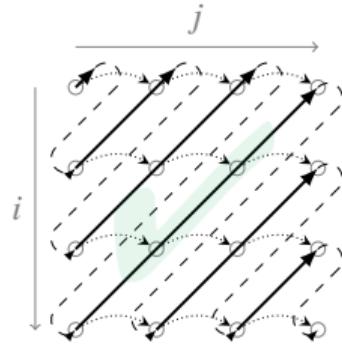
```
for(int i = 0; i < N; i++)
    for(int j = 1; j < M; j++)
        A[i, j] += A[i, j-1]; // S[i][j]
```



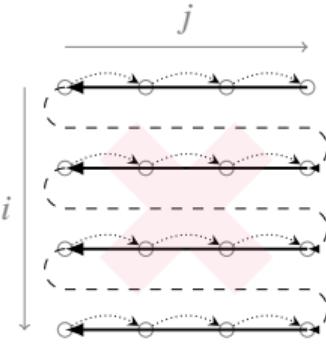
a $\theta(S[i, j]) = (i, j)$



b $\theta(S[i, j]) = (j, i)$



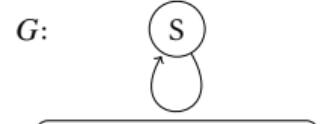
c $\theta(S[i, j]) = (i + j, j)$



d $\theta(S[i, j]) = (i, -j)$

Legality check

G:



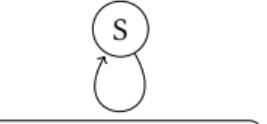
$$S[i, j - 1] \mapsto S[i, j]$$

$$\theta(S[i, j]) = (i + j, j)$$

$$(i + j - 1, j - 1) \mapsto (i + j, j)$$

$$\begin{aligned} \delta &= (i + j, j) - (i + j - 1, j - 1) \\ &= (1, 1) >_{\text{LEX}} \vec{0} \end{aligned}$$

G:



$$S[i, j - 1] \mapsto S[i, j]$$

$$\theta(S[i, j]) = (i, -j)$$

$$(i, -j + 1) \mapsto (i, -j)$$

$$\begin{aligned} \delta &= (i, -j) - (i, -j + 1) \\ &= (0, -1) \not>_{\text{LEX}} \vec{0} \end{aligned}$$

Symbolic representation

$$\{(i, j) : 1 \leq i \leq N \wedge 1 \leq j \leq i\}$$

$$1 \leq i$$

$$i \leq N$$

$$1 \leq j$$

$$j \leq i$$

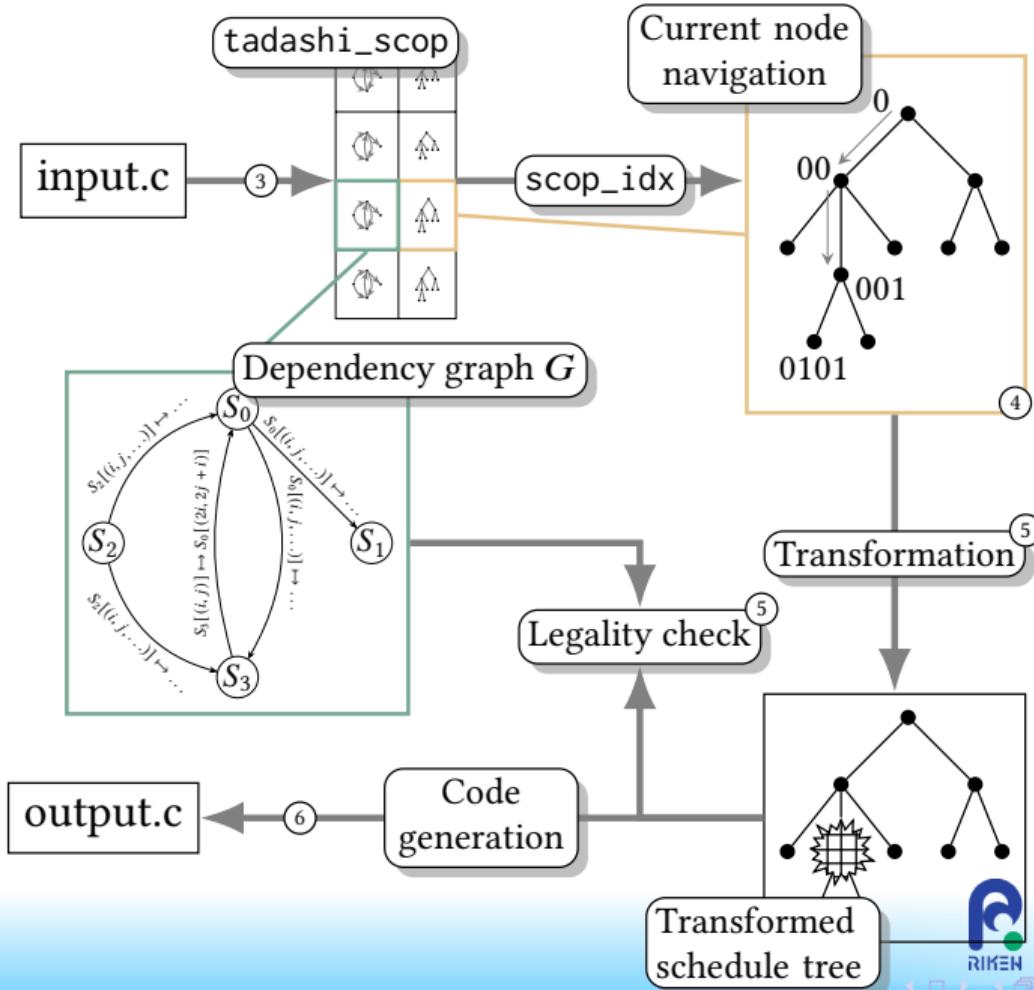
$$0 \leq i - 1$$

$$0 \leq N - i$$

$$0 \leq j - 1$$

$$0 \leq i - j$$

N	i	j	const
0	1	0	-1
1	-1	0	0
0	0	1	-1
0	1	-1	0



That's all folks!

- ▶ i. <https://vatai.github.io/>
- ▶ ii. <https://github.com/vatai/tadashi/>
- ▶ iii. <https://arxiv.org/abs/2410.03210>

