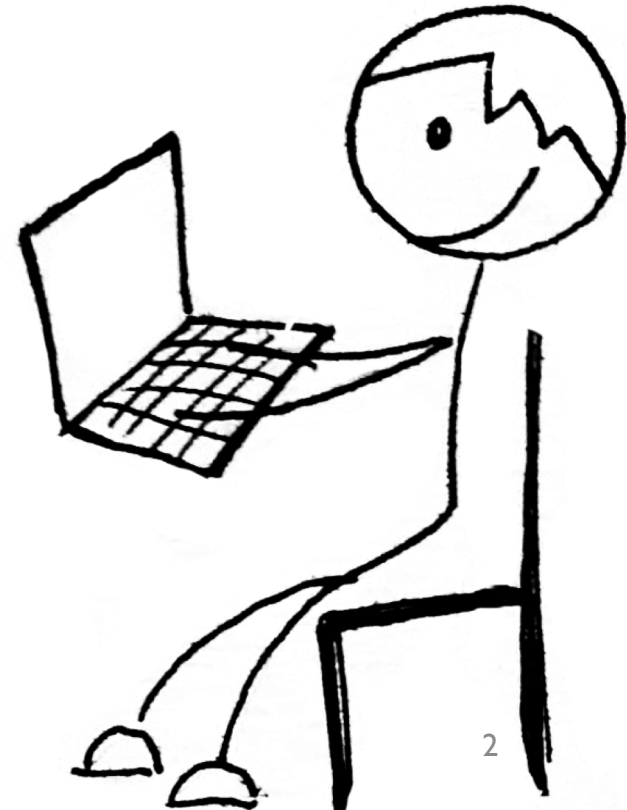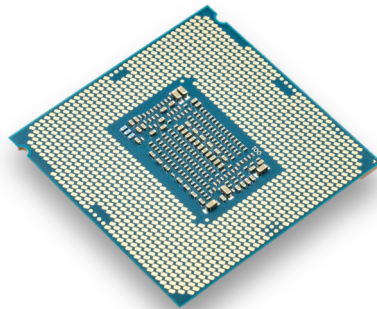# Toward Self-Evolving Compilers

Phitchaya Mangpo Phothilimthana

Google Brain
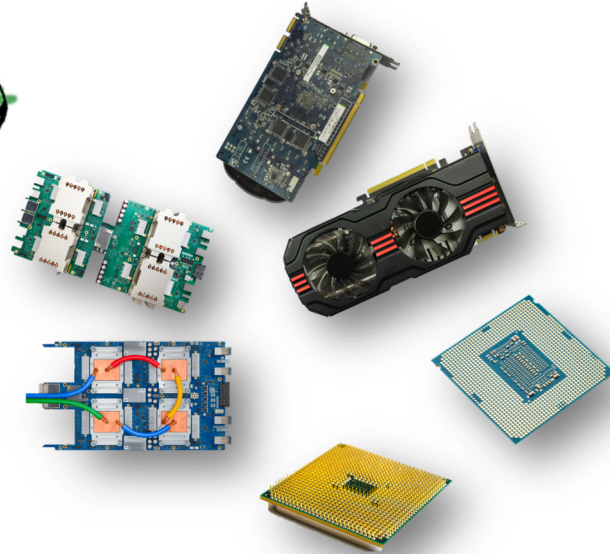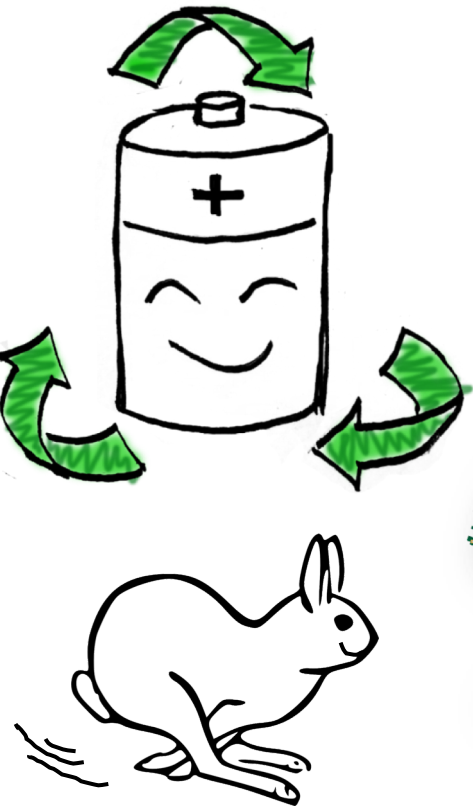
*based on the work of **many** people*

unusual ISA

heterogeneity

no register

restricted computations
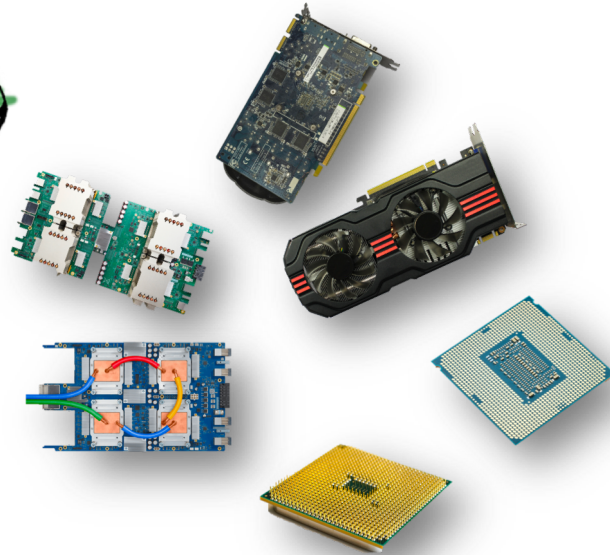
limited resources

new memory hierarchy

distributed memory

3

# unusual ISA

## heterogeneity

### no register

restricted computations

## limited resources

### new memory hierarchy

## distributed memory

good programming model
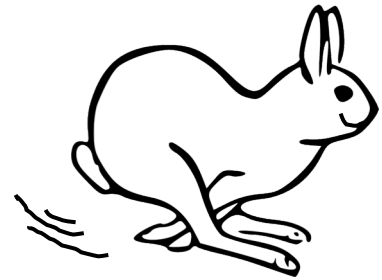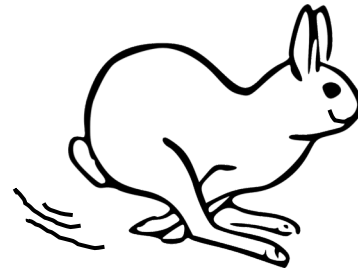
**unusual ISA**

heterogeneity

no register

restricted computations

limited resources

new memory hierarchy

distributed memory

good programming model

5

**unusual ISA**

**heterogeneity**

**no register**

**restricted computations**

**limited resources**

new memory hierarchy

**distributed memory**

good programming model

**Synthesis-aided compilation**

**unusual ISA**

**heterogeneity**

**no register**

**restricted computations**

**imited resources**

**new memory hierarchy**

**distributed memory**

# Classical

input program $p$

$p'$

$p''$

...



correct

all programs

**Pros:** fast to compile

**Cons:** miss efficient programs

# Synthesis-aided

# Classical

input program $p$

$p'$

$p''$

...

correct

all programs

**Pros:** fast to compile

**Cons:** miss efficient programs

# Synthesis-aided

Define search space $P$
Find $p \in P$

program $p$ $\equiv$ input program $p_{spec}$

correct

Search space

all programs

**Pros:** find provably optimal program

**Cons:** slow to compile

**Swizzle Inventor**

semi-automatic
bypass rewrite rules & heuristics



**Superoptimization**

fully-automatic
bypass rewrite rules & heuristics



**Auto** 

fully-automatic
bypass heuristics

# Swizzle
# Inventor

semi-automatic
bypass rewrite rules & heuristics

**GreenThumb**
Superoptimization

fully-automatic
bypass rewrite rules & heuristics

**AutoX**

fully-automatic
bypass heuristics

# Swizzle

**non-trivial movement** of data or
**non-trivial mapping** of computations
to **hardware resources** and **loop iterations**

for dramatic performance improvement

# Register Cache: Stencil



global memory

load

shared memory

registers

output

output

t0   t1   t2   t3

t0   t1   t2   t3

*[Ben-Sasson et al. ICS' 16]*

# Register Cache: Stencil

global memory

shared memory

registers

output

output

# Register Cache: Stencil

*[Ben-Sasson et al. ICS' 16]*

# Register Cache: Stencil



global memory

load

shared memory    registers

output

output

t0  t1  t2  t3

t0  t1  t2  t3

**In each iteration**

```
__shfl_sync(mask, rc[idx],
            recv_from)
```

**rc**

| a | b |

**idx:**
(tid >= k)? 0 : 1

**recv_from:**
(tid + k) % warpSize

*[Ben-Sasson et al. ICS' 16]*

16

# **Swizzle** Inventor

Helps programmers implement swizzle programs by:

- letting them **write program sketches that omit swizzles**

- **automatically synthesizing swizzles** to complete the programs

# Stencil: Program Sketch

```
rc = load(input, warpOffset,
        /* slice */ 1,
        /* iterations */ 2);

int out = 0;
for(int k = 0; k < 3; k++) {
    int tmp = magic_get(rc);
    out += tmp;
}
```

**output**

| 0 | 1 | 2 | 3 |

t0    t1    t2    t3

```
output[tid] = out;
```

# Stencil: Program Sketch



global memory

**input** | 0 | 1 | 2 | 3 | 4 | 5 |

registers

**rc** | 0 | 4 | 1 | 5 | 2 | | 3 | |

**output** | 0 | 1 | 2 | 3 |

t0  t1  t2  t3

```
rc = load(input, warpOffset,
    /* slice */ 1,
    /* iterations */ 2);

int out = 0;
for(int k = 0; k < 3; k++) {
    int tmp = magic_get(rc);
    out += tmp;
}

output[tid] = out;
```

# Stencil: Program Sketch

global memory

**input**

| 0 | 1 | 2 | 3 | 4 | 5 |

registers

**rc**

| 0 | 4 | 1 | 5 | 2 | | 3 | |

**output**

| 0 | 1 | 2 | 3 |

t0  t1  t2  t3

```
rc = load(input, warpOffset,
        /* slice */ 1,
        /* iterations */ 2);

int out = 0;
for(int k = 0; k < 3; k++) {
    int tmp = magic_get(rc);
    out += tmp;
}

output[tid] = out;
```

# Stencil: Program Sketch

```
int tmp = magic_get(rc); -->

// Choose which input data to send
int idx = ?sw_part(2, tid, k);

// Choose which thread to read from
int recv_from =
?sw_xform(tid, warpSize, k);

// Perform intra-warp shuffle
int tmp = __shfl_sync(FULL_MASK, rc[idx], recv_from);
```

# Transformation Swizzle Hole

**?sw_xform** hole defines the search space that contains **grouping** permutations of **fanning** followed by **rotation**.

**rotation**

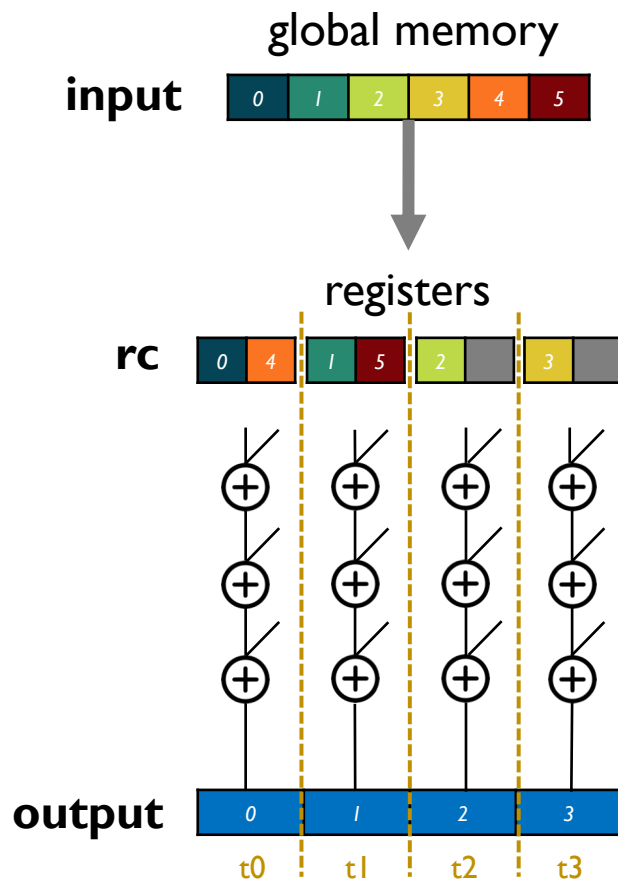$$rot(\mathbf{i}) = (\mathbf{i}+2) \bmod \mathbf{8}$$

$x[i]$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7

$y[rot(i)] = x[i]$ | 6 | 7 | 0 | 1 | 2 | 3 | 4 | 5

co-prime **fanning**

fan size

$$fan(\mathbf{i}) = (\mathbf{3}*\mathbf{i}) \bmod \mathbf{8}$$

$x[i]$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7

$y[fan(i)] = x[i]$ | 0 | 3 | 6 | 1 | 4 | 7 | 2 | 5

fan size

**grouping**

$$group(\mathbf{4}, fan)(\mathbf{i}) = \lfloor \mathbf{i}/\mathbf{4} \rfloor * \mathbf{4} + ((\mathbf{3}*(\mathbf{i} \bmod \mathbf{4})) \bmod \mathbf{4})$$

$x[i]$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7

$y[group(4, fan)(i)] = x[i]$ | 0 | 3 | 2 | 1 | 4 | 7 | 6 | 5

$gs$

# Correctness Condition

**Spec: sequential program**

```
void spec(
    const float *x,
    float *y, int n) {

  for(int i = 0; i < n; i++) {
    int out = 0;
    for(int k = 0; k < 3; k++)
      out += x[i+k];
    y[i] = out;
  }
}
```

$$\exists h \, \forall x \,.\, spec(x, y, n)$$
$$\wedge \, sketch(h)(x, y', n)$$
$$\wedge \, y = y'$$

**Sketch: CUDA sketch**

```
__global__ void sketch(
    const float *x,
    float *y, int n) {

  rc = load(x, warpOffset, 1, 2);

  int out = 0;
  for(int k = 0; k < 3; k++) {
    int tmp = magic_get(rc);
    out += tmp;
  }

  y[tid] = out;
}
```

# **Inventiveness**:
## Can Swizzle Inventor invent new optimizations?

# Finite Field Multiplication

```
// Create ans0, ans1, ans2, ans3
acc ans0 = create_accumulator(0, identity, ^, &);
...

for(int k = 0; k < 32; k++) {
  int a0 = __shfl_sync(mask, rA[?sw_part(2,tid,k)],
                       ?sw_xform(tid,32,k));
  int a1 = __shfl_sync(mask, rA[?sw_part(2,tid,k)],
                       ?sw_xform(tid,32,k));
  int b0 = __shfl_sync(mask, rB[?sw_part(2,tid,k)],
                       ?sw_xform(tid,32,k));
  int b1 = __shfl_sync(mask, rB[?sw_part(2,tid,k)],
                       ?sw_xform(tid,32,k));

  // Update ans0
  accumulate(ans0, [a0,b0], ?sw_cond(tid,k));
  accumulate(ans0, [a0,b1], ?sw_cond(tid,k));
  accumulate(ans0, [a1,b0], ?sw_cond(tid,k));
  accumulate(ans0, [a1,b1], ?sw_cond(tid,k));

  // Update ans1, ans2, ans3
  ...
}
```
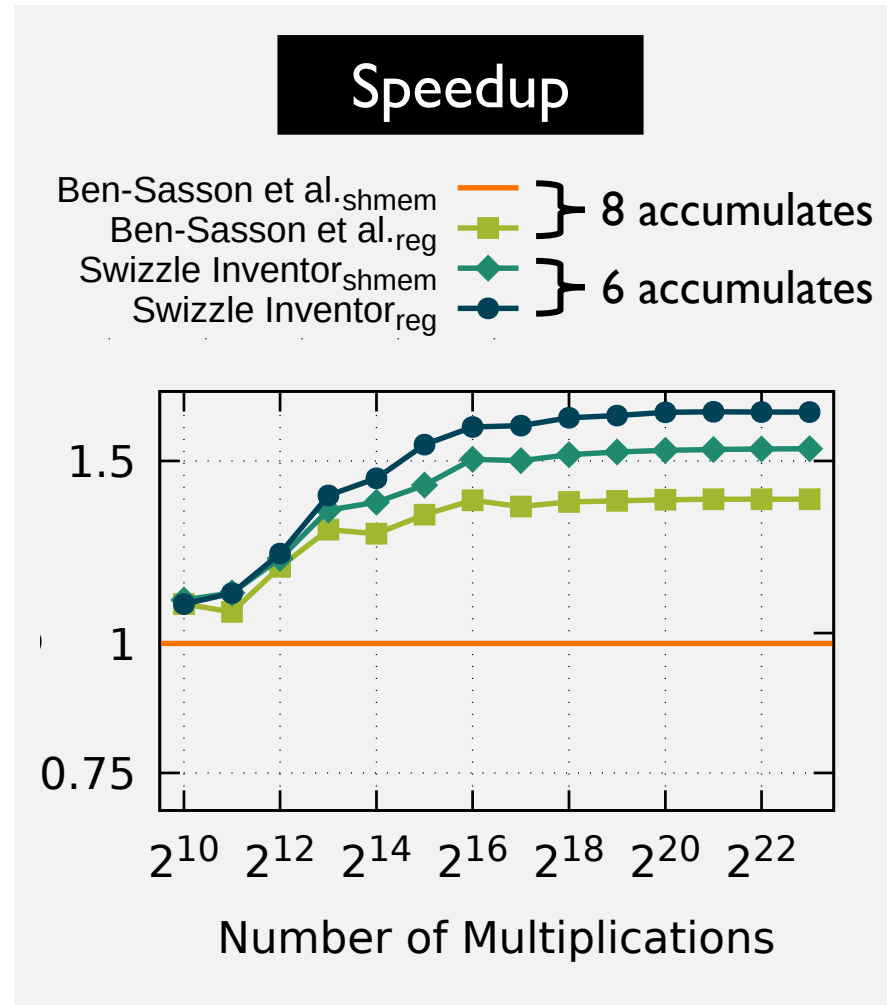


Speedup

Ben-Sasson et al.$_{shmem}$
Ben-Sasson et al.$_{reg}$ } 8 accumulates
Swizzle Inventor$_{shmem}$
Swizzle Inventor$_{reg}$ } 6 accumulates

Number of Multiplications

*Baseline: [Ben-Sasson et al. ICS' 16]*

# Matrix Transposition



global memory

load

New algorithm!

registers

Swizzle Inventor synthesizes in seconds!

Search space = ~$10^{23}$

column shuffle

row shuffle

row shuffle

column shuffle

column rotate

row shuffle

*[Catanzaro et al. PPoPP '14]*

# Swizzle Inventor

semi-automatic
bypass rewrite rules & heuristics

## Green Thumb
### Superoptimization

fully-automatic
bypass rewrite rules & heuristics

## AutoX

fully-automatic
bypass heuristics

# Swizzle
# Inventor

**Green Thumb**

### Superoptimization

**fully-automatic
bypass rewrite rules & heuristics**

# AutoX

fully-automatic
bypass heuristics

29

## gcc –O3

**ARM**
register-based ISA

```
cmp    r1, #0
mov    r3, r1, asr #31
add    r2, r1, #7
mov    r3, r3, lsr #29
movge  r2, r1
ldrb   r0, [r0, r2, asr #3]
add    r1, r1, r3
and    r1, r1, #7
sub    r3, r1, r3
asr    r1, r0, r3
and    r0, r0, #1
```



## 82% speedup

```
asr    r3, r1, #2
add    r2, r1, r3, lsr #29
ldrb   r0, [r0, r2, asr #3]
and    r3, r2, #248
sub    r3, r1, r3
asr    r1, r0, r3
and    r0, r1, #1
```

**GreenArrays**
stack-based ISA

## Expert's

```
push over - push and
pop pop and over
0xffff or and or
```

Precondition: top 3 elements in
the stack are <= 0xffff



## 2.5X speedup

```
dup push or and pop or
```

**GOAL:** develop a **search technique** that can synthesize optimal programs **faster** and more **consistently**.

# We developed …

**1** **Lens**
enumerative
search algorithm

# **Lens** Enumerative Search Algorithm

# We developed …



inst1 … …   $p_{pre}$
inst2 … …
inst3 … …
inst4 … …   $p_{spec}$
inst5 … …
inst6 … …
inst7 … …   $p_{post}$

**1** **Lens**
enumerative
search algorithm

OR
symbolic search
stochastic search

**2** **Context-aware
window decomposition**

**3** **Cooperative
search instances**

**shared
data**

# Runtime Speedup

Runtime speedup over **gcc –O3** on an actual **ARM Cortex-A9**

**Benchmarks** Hacker's Delight, WiBench (wireless system kernel benchmarks), MiBench (embedded system kernel benchmarks)

| Program | Search time (s) | gcc –O3 length | Output length | Runtime speedup on ARM Cortex-A9 |
|---|---|---|---|---|
| p18 | 9 | 7 | 4 | **2.11** |
| p21 | 1139 | 6 | 5 | **1.81** |
| p23 | 665 | 18 | 16 | **1.48** |
| p24 | 151 | 7 | 4 | **2.75** |
| p25 | 2 | 11 | 1 | **17.80** |
| WB-txrate5a | 32 | 9 | 8 | **1.31** |
| WB-txrate5b | 66 | 8 | 7 | **1.29** |
| MB-bitarray | 612 | 10 | 6 | **1.82** |
| MB-bitshift | 5 | 9 | 8 | **1.11** |
| MB-bitcnt | 645 | 27 | 19 | **1.33** |
| MB-susan-391 | 32 | 30 | 21 | **1.26** |

# *GreenThumb* Framework

Provide cooperative search strategy.

Enable rapid retargeting of the superoptimizer to a new ISA.

github.com/mangpo/greenthumb

**Supported ISAs:** GreenArrays, ARM, subset of LLVM

**LinkiTools' S10** (https://linki.tools/s10.html): RISC-V

# Swizzle Inventor

Green Thumb
## Superoptimization

**fully-automatic
bypass rewrite rules & heuristics**

AutoX

fully-automatic
bypass heuristics

37

# Swizzle Inventor

semi-automatic
bypass rewrite rules & heuristics

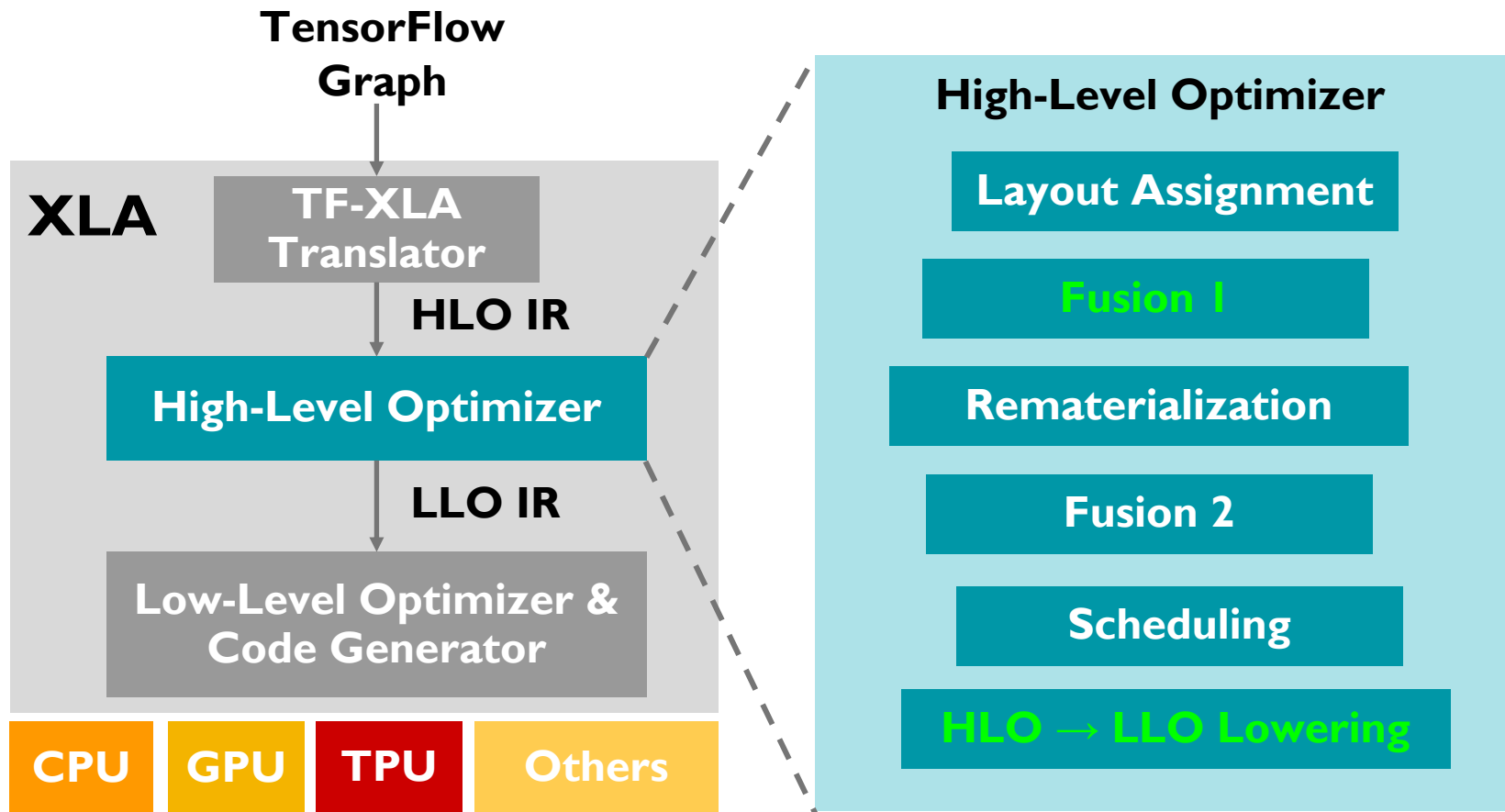# GreenThumb
## Superoptimization

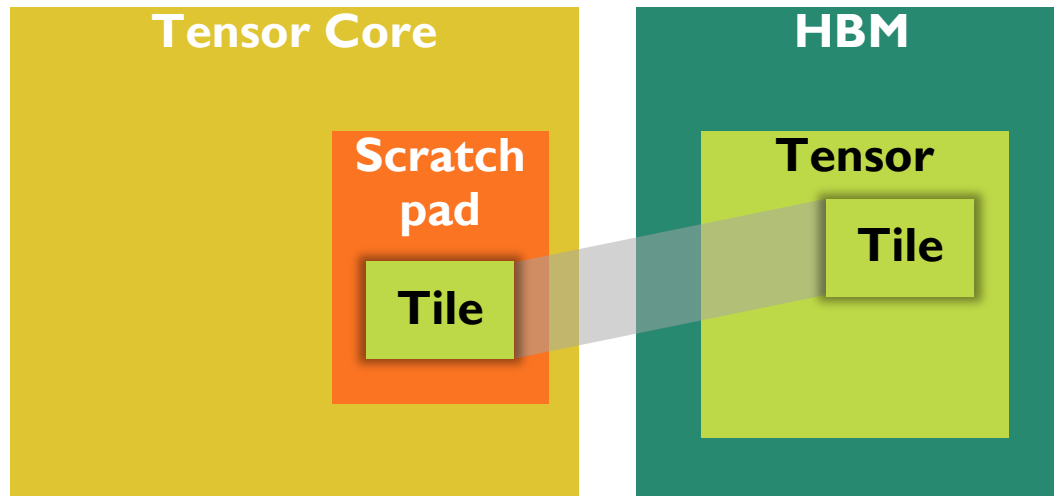fully-automatic
bypass rewrite rules & heuristics

# AutoX

**fully-automatic
bypass heuristics**

# XLA: Accelerated Linear Algebra

# Tile Size in Lowering Pass



TPUs process one (fused) tensor op at a time.

- Entire tensors don't fit in scratchpad.
- To process one tile of an output tensor, copy input tiles into scratchpad.
- Store intermediates in scratchpad.

# Autotuning Tile Size

## Autotuner
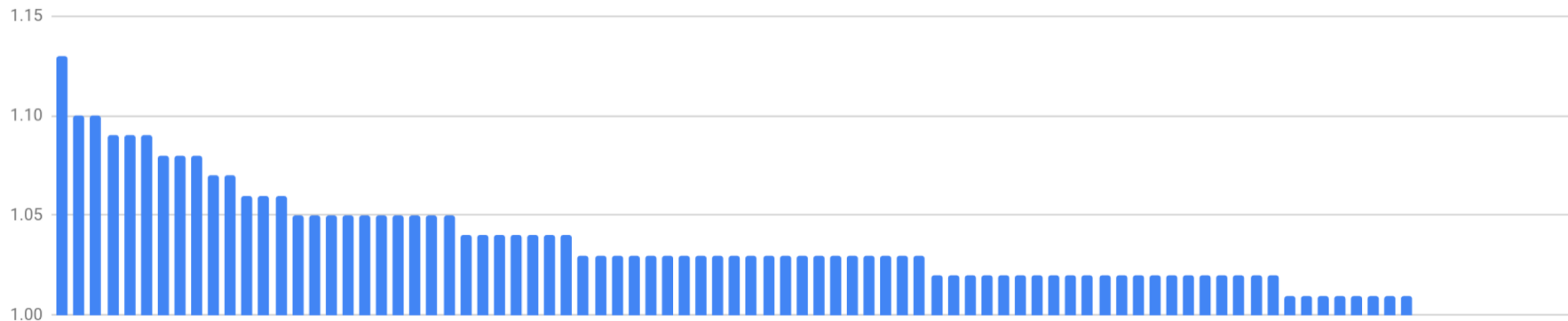
- Exhaustive search (100-1M choices per program)
- Evaluate by running on real hardware.
- Fast mode: tune subset of candidates.
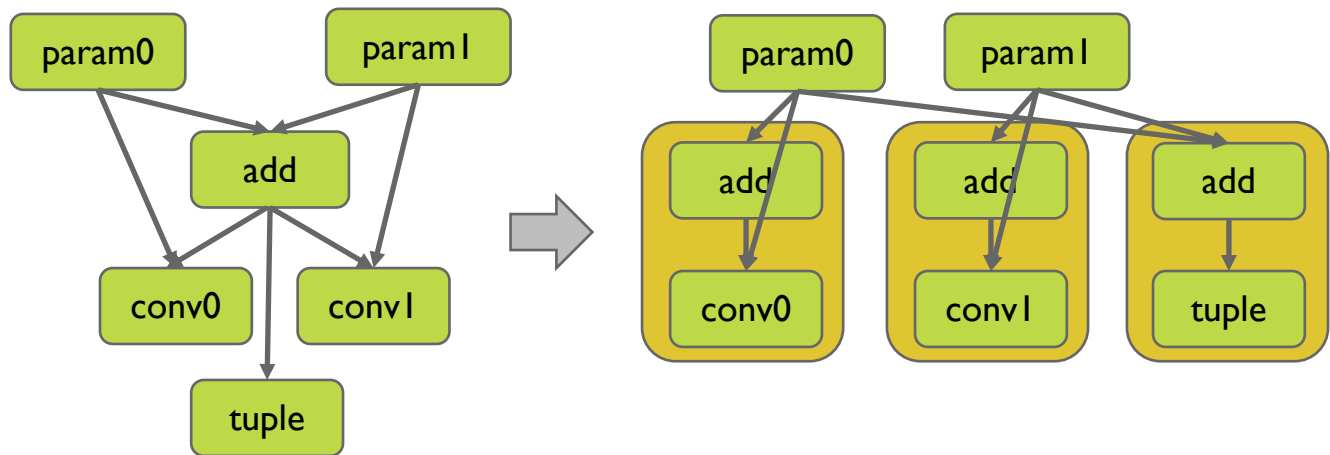
## Result:

24 benchmarks gain >= 5% improvement

Speedup compared to default

# Fusion Decisions

Example:



Fusion configuration:
- fuse or do-not-fuse per **node**
- If a node is marked fuse, it is fused into all its consumers.
- 100 - 100,000 nodes per graph

# Autotuning Fusion Decisions

## Autotuner

- Partition graph into <=1000-node clusters.

- Run simulated annealing in each cluster.

- Start the search from the default or random configuration.

## Highlighted results

- 9% speedup on search ranking inference model

- 13% speedup on TPU compute time of ad recommendation training model

- 6-9% speedup on ML Perf inference models

**Swizzle Inventor**

semi-automatic
bypass rewrite rules & heuristics

**GreenThumb Superoptimization**

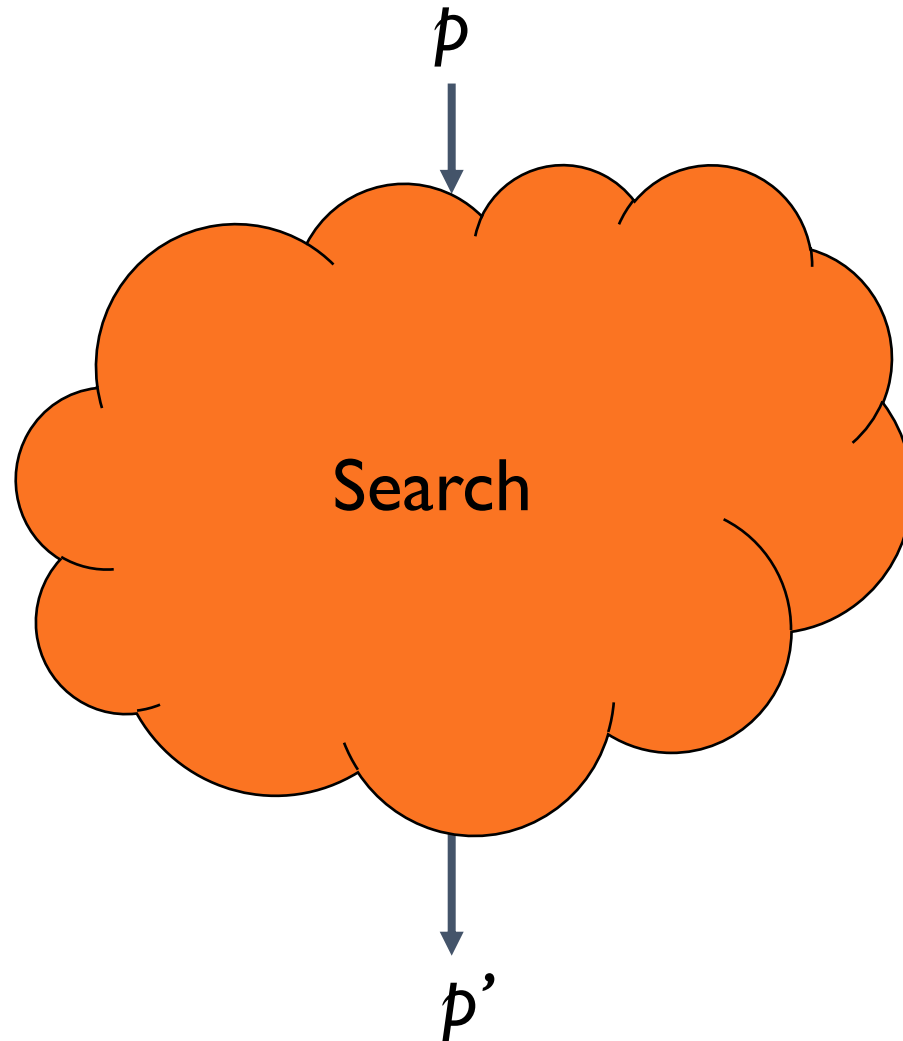fully-automatic
bypass rewrite rules & heuristics

**AutoX**

fully-automatic
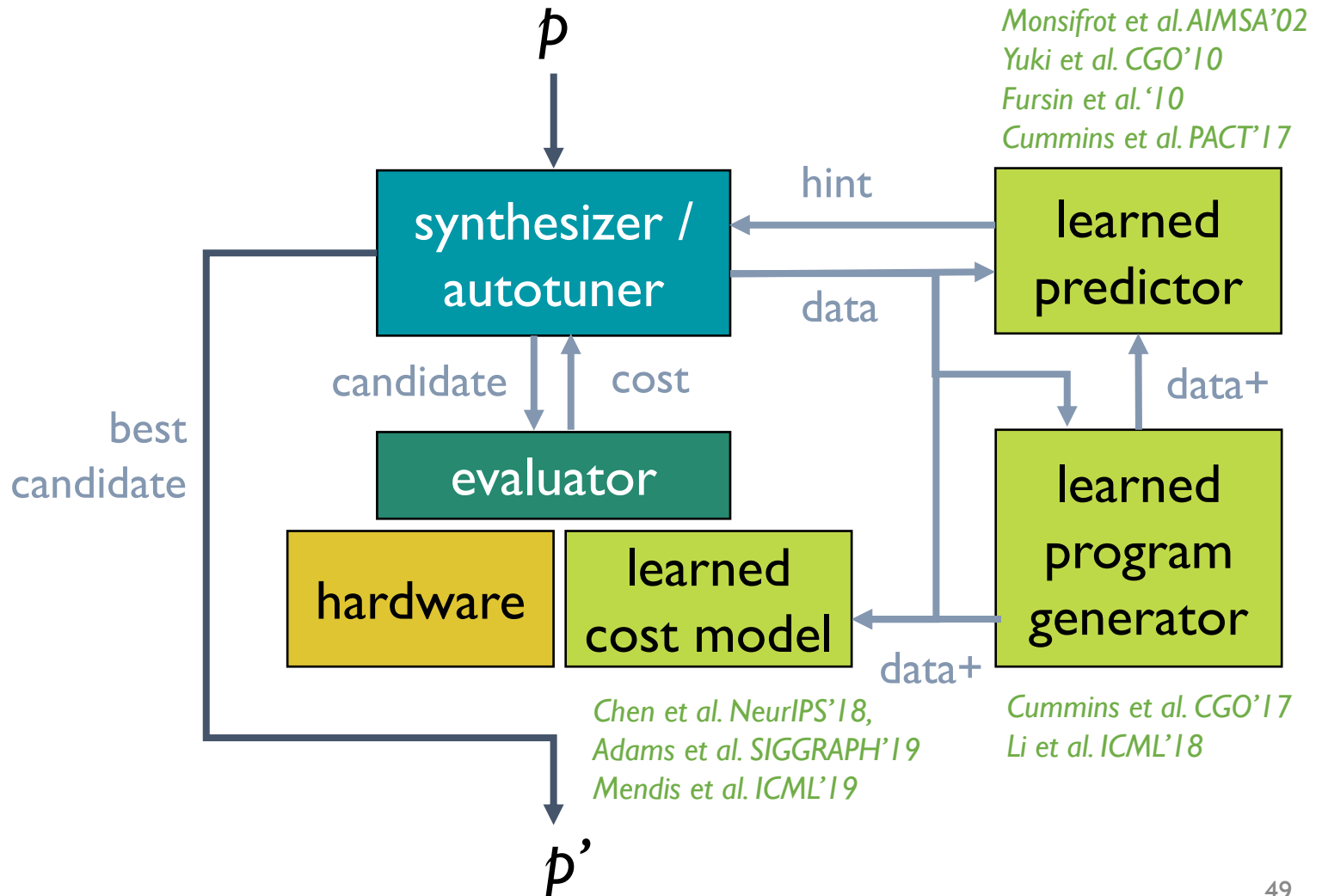bypass heuristics

# Problems of Search-Based Compilers

slow

do not learn from past experience

# Toward Self-Evolving Compilers

$p$

Search

$p'$

# Toward Self-Evolving Compilers



$p$

*Monsifrot et al. AIMSA'02*
*Yuki et al. CGO'10*
*Fursin et al.'10*
*Cummins et al. PACT'17*

synthesizer / autotuner

hint

learned predictor

data

best candidate

candidate

cost

data+

evaluator

hardware

learned cost model

learned program generator

data+

*Chen et al. NeurIPS'18,*
*Adams et al. SIGGRAPH'19*
*Mendis et al. ICML'19*

*Cummins et al. CGO'17*
*Li et al. ICML'18*

$p'$

49

# Toward Self-Evolving Compilers

**Is it possible** to build

a self-evolving compiler that is

**as fast as a heuristics-base** compiler

and **produces better code**?

**Swizzle Inventor**

**GreenThumb**

**Auto**

Archibald Samuel Elliott
An Wang
Abhinav Jangda
Bastian Hagedorn
Henrik Barthels
Samuel J. Kaufman
Vinod Grover
Emina Torlak
Rastislav Bodik

Aditya Thakur
Rastislav Bodik
Dinakar Dhurjati

Mike Burrows
Bjarke Roune
Amit Sabne
Sam J. Kaufman
Cambridge Yang
Dimitrios Vytiniotis
Dominik Grewe

# Collaborators