

Deep Learning for Compilers, Fun, & Profit

Chris Cummins
<https://chriscummins.cc>
cummins@fb.com

overview

- 1. What I used to do**
- 2. What I now do**
- 3. What you should do**

overview

- 1. What I used to do**
- 2. What I now do**
- 3. What you should do**



I'm not qualified to give advice

Every PhD is different

Sample size = 1

My PhD credentials:

I finished

I achieved my goals

I had fun

Good reasons to do a PhD

-  You don't want a real job
-  You don't want to grow up
-  You want money

What is a PhD

- "on-the-job training for being a scientist"
- pure transferable skills
- a thesis

What is a thesis

- Thesis = 3 papers 
- 1 paper  = 1 idea 
- ergo, PhD = 3 ideas

What is a PhD NOT

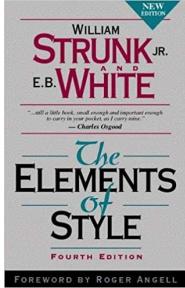
- an excuse to be a loner
- validation of your intellect
- a grind

What you should do

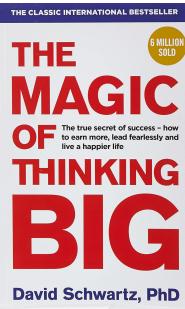
- 1. Figure out what success looks like:**
 - to your viva panel**
 - to yourself**
- 2. Don't think of a PhD as a job**
- 3. Travel!**

Reading list

1.



2.



2. 1500+ research papers



What you should do: industry transition

- 1. Take internships**
- 2. Start practicing interviews early & often**
- 3. Don't get 1 job offer, aim for 3**

overview

- 1. What I used to do**
- 2. What I now do**
- 3. What you should do**

overview

- 1. What I used to do**
- 2. What I now do**
- 3. What you should do**

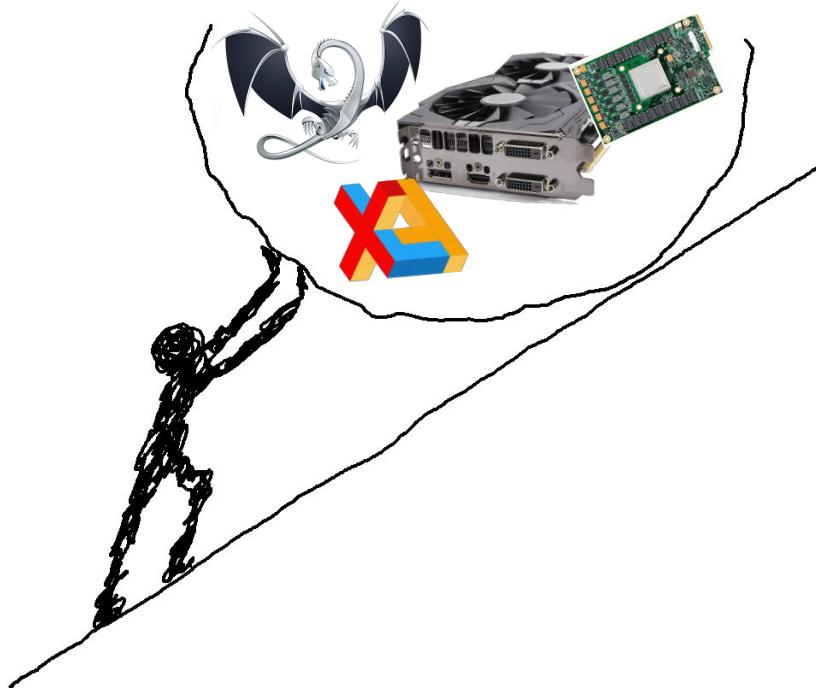
"machine learning for compilers for machine learning"

A Venn diagram consisting of two overlapping circles. The left circle is pink and labeled "Compilers". The right circle is light blue and labeled "Machine Learning". The intersection of the two circles contains a small circular inset featuring a smiling man giving a thumbs-up. The entire diagram is set against a white background.

Compilers

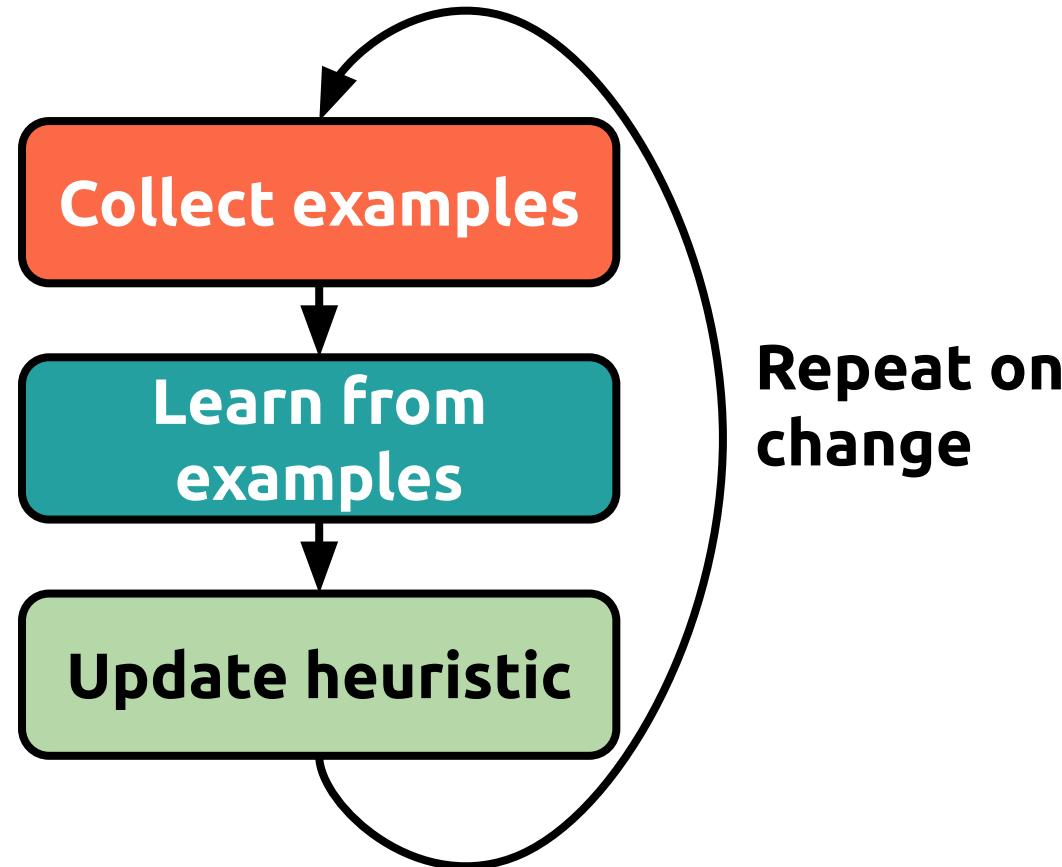
**Machine
Learning**

Building compilers... a job for life



- Compiler / HW keeps changing
- 1000s of variables
- NP-hard or worse
- Bad heuristics
- Wasted energy
- Widening performance gap

"Build an optimizing compiler, your code will be fast for a day.
Teach a compiler to optimize ... "



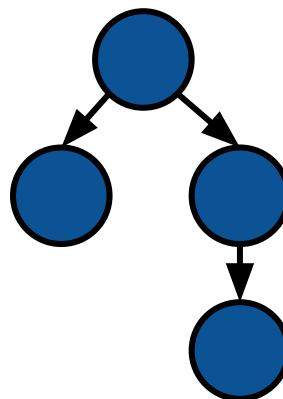
Summarize the program

Program

```
void LinearAlgebraOp<InputScalar,  
OutputScalar>::AnalyzeInputs(  
    OpKernelContext* context, TensorInputs* inputs,  
    TensorShapes* input_matrix_shapes, TensorShape*  
batch_shape) {  
    int input_rank = -1;  
    for (int i = 0; i < NumMatrixInputs(context); ++i) {  
        const Tensor& in = context->input(i);  
        if (i == 0) {  
            input_rank = in.dims();  
            OP_REQUIRES(  
                context, input_rank >= 2,  
                errors::InvalidArgument(  
                    "Input tensor ", i,  
                    " must have rank >= 2"));  
    }
```



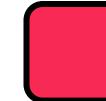
IR (CFG, DFG, AST,...)



Features



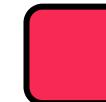
#. instructions



loop nest level



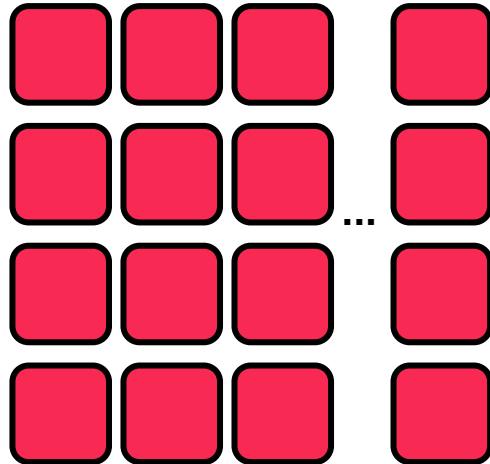
arithmetic intensity



trip counts

Collect examples

Features

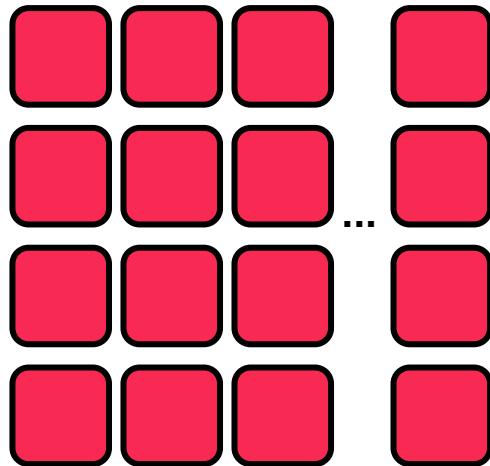


Best Param

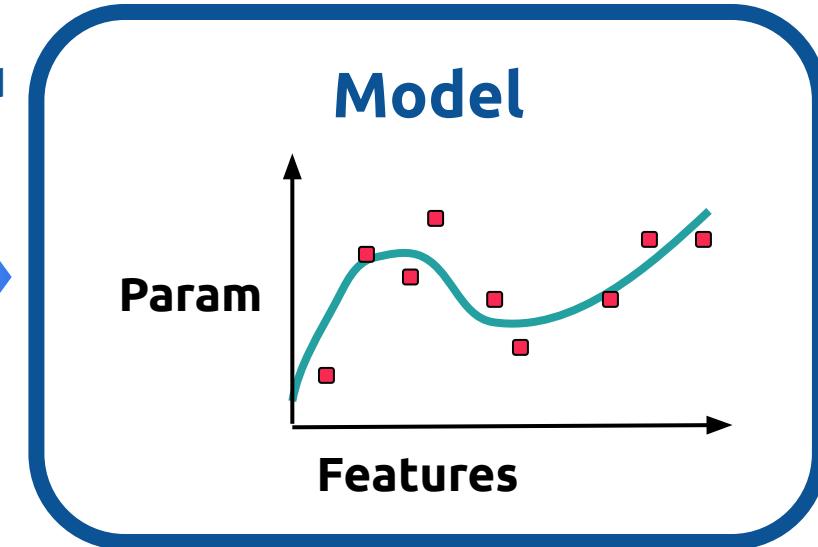


Learn from examples

Features



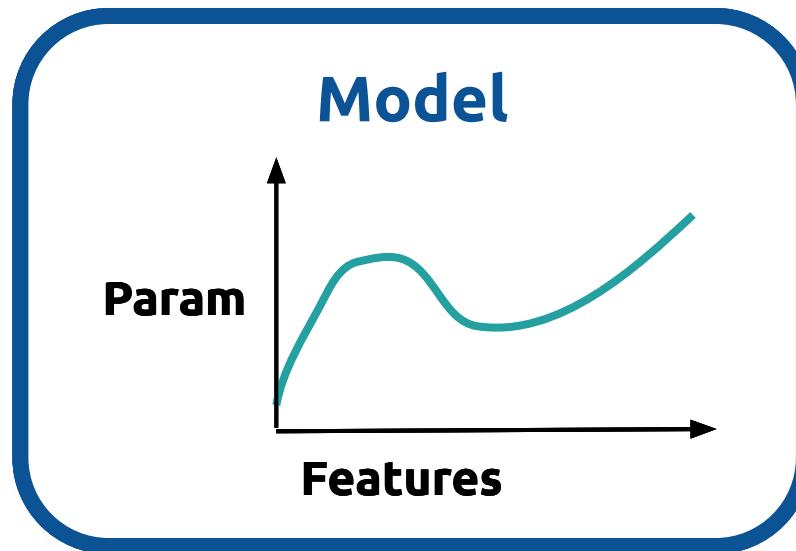
Supervised
Machine
Learner



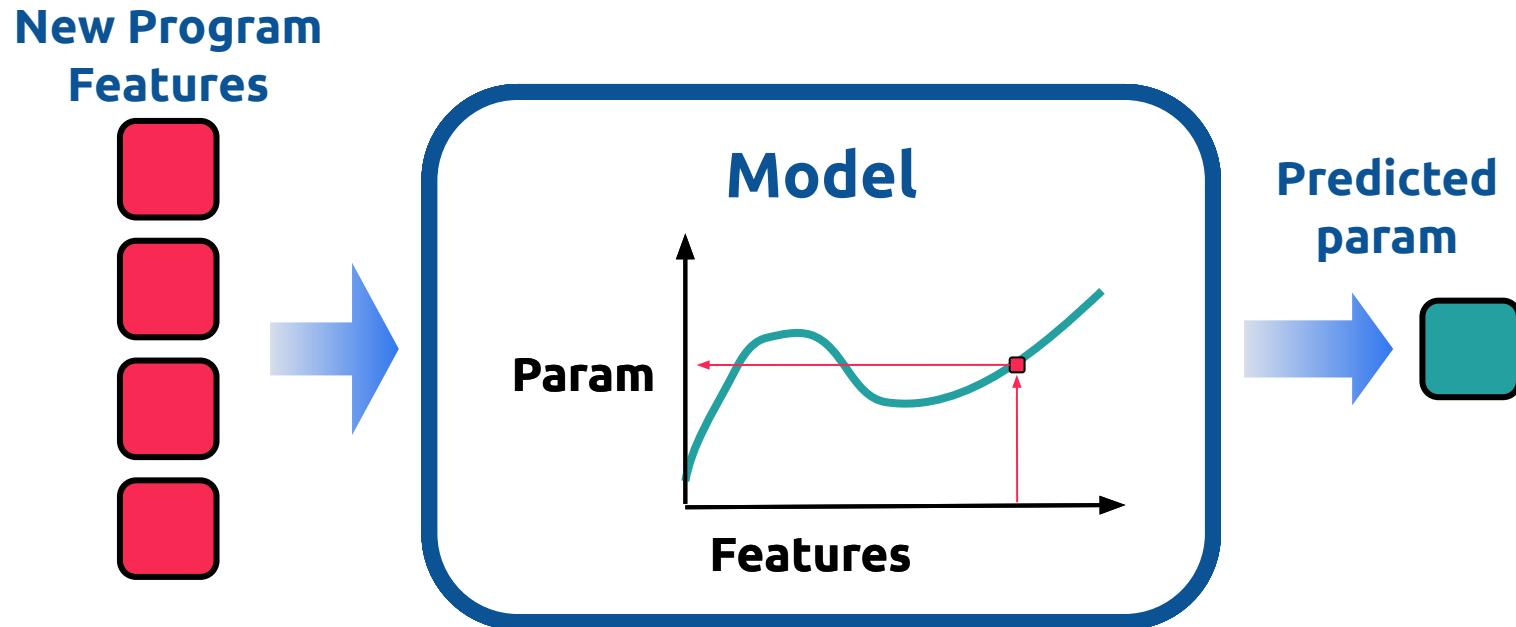
Best Param



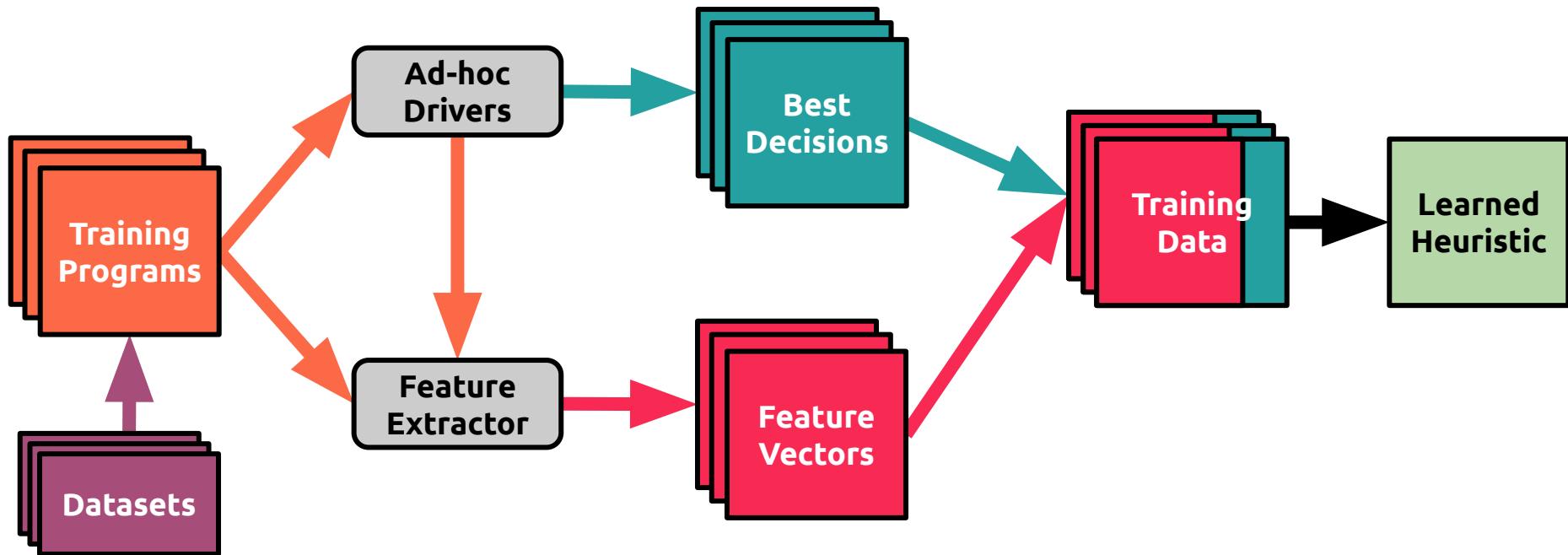
The model is the heuristic



The model is the heuristic

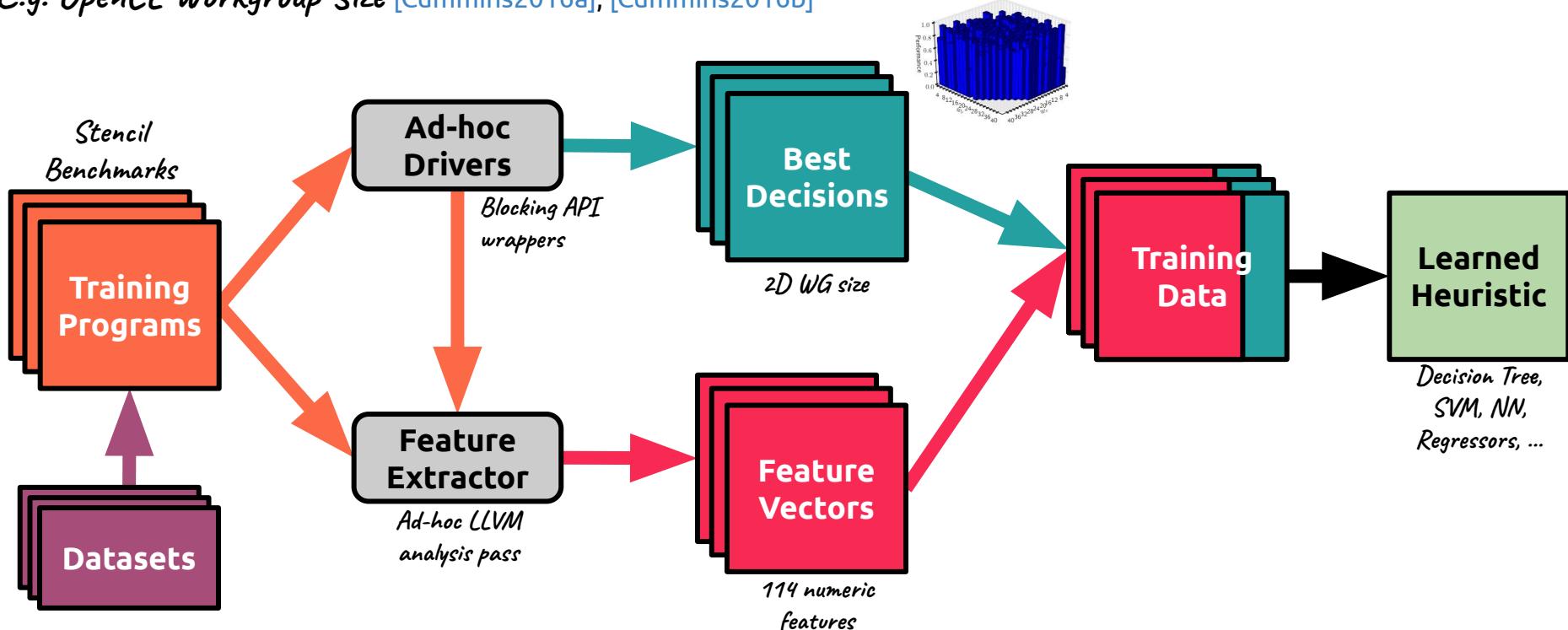


Machine learning for compilers



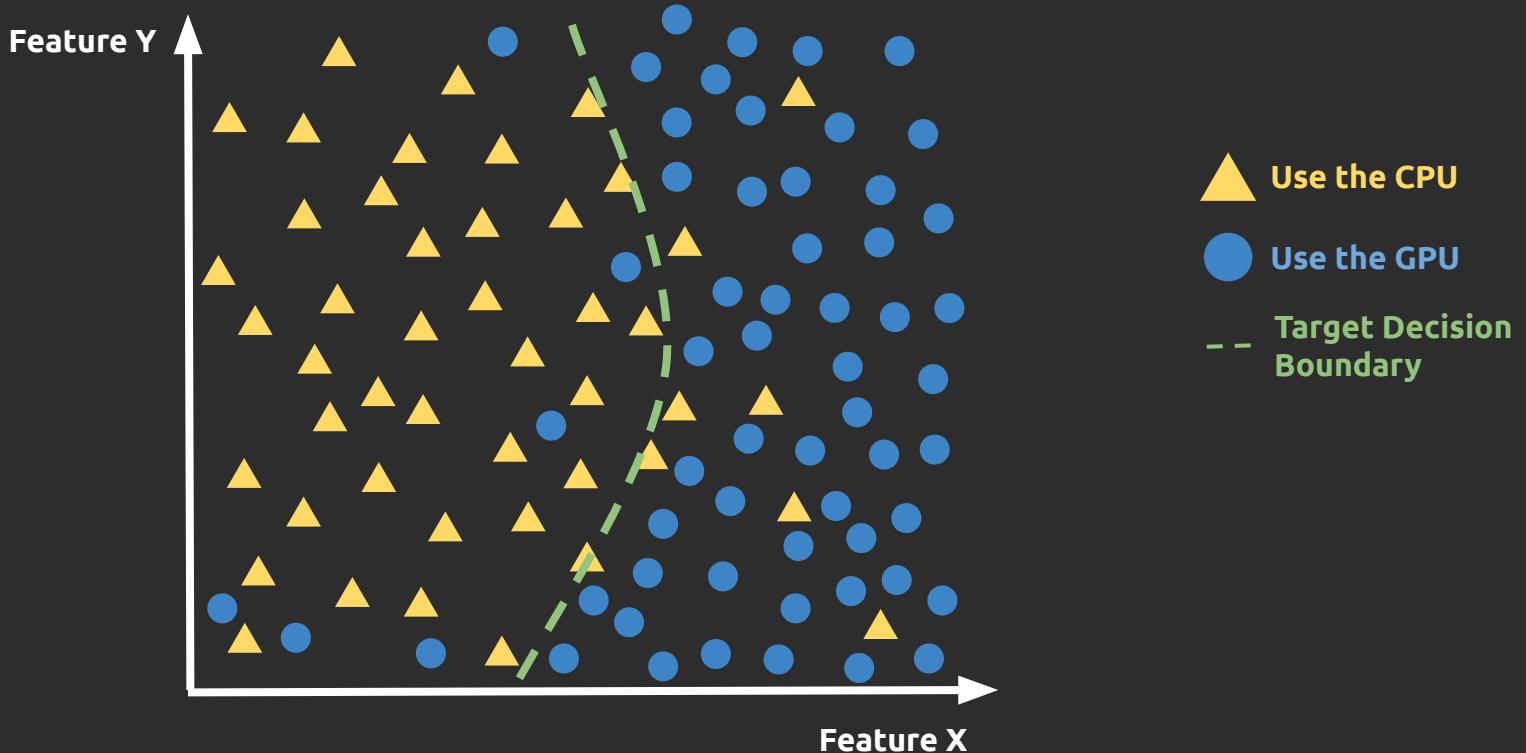
Machine learning for compilers

E.g. OpenCL Workgroup Size [Cummins2016a], [Cummins2016b]

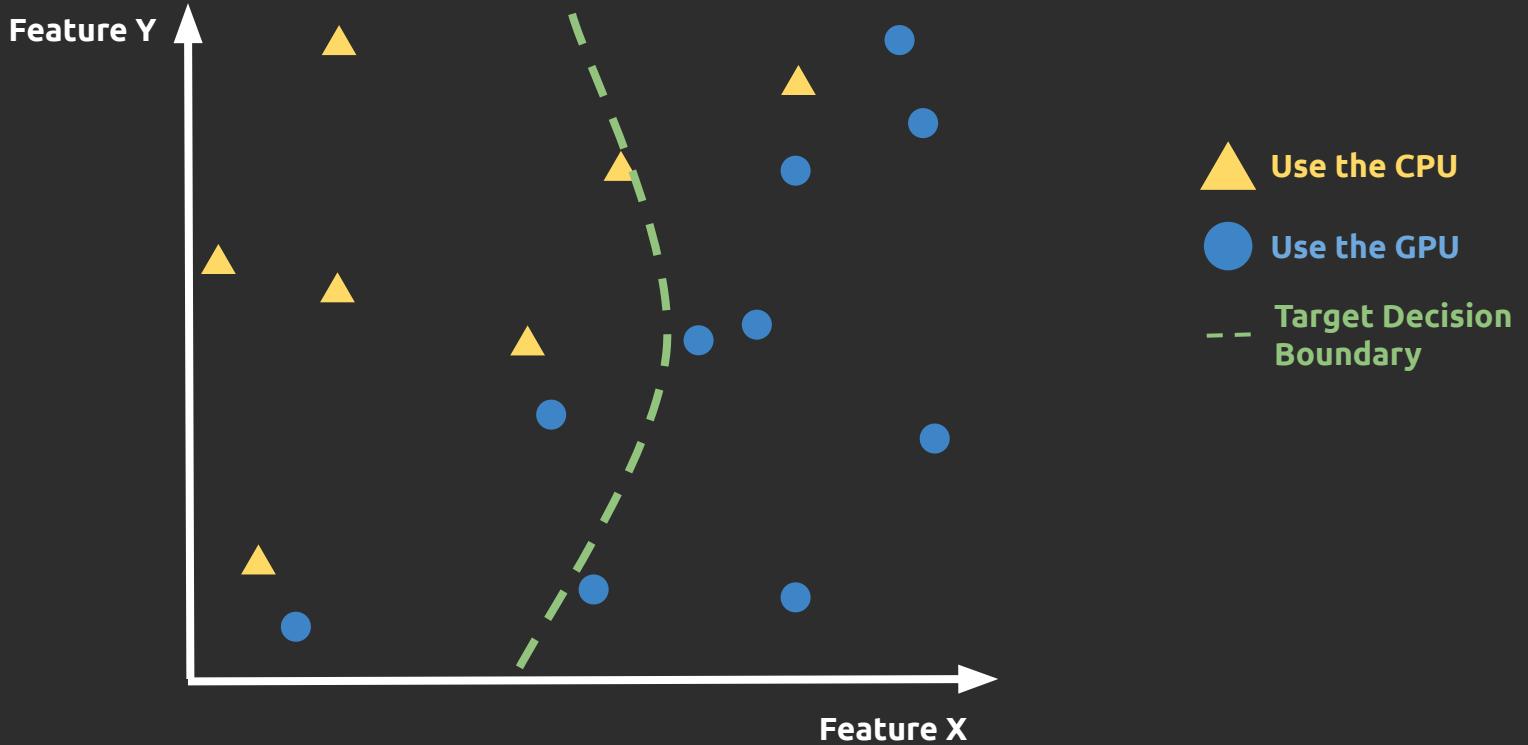


*3.34x better than default heuristic
18% better than domain experts*

What we want



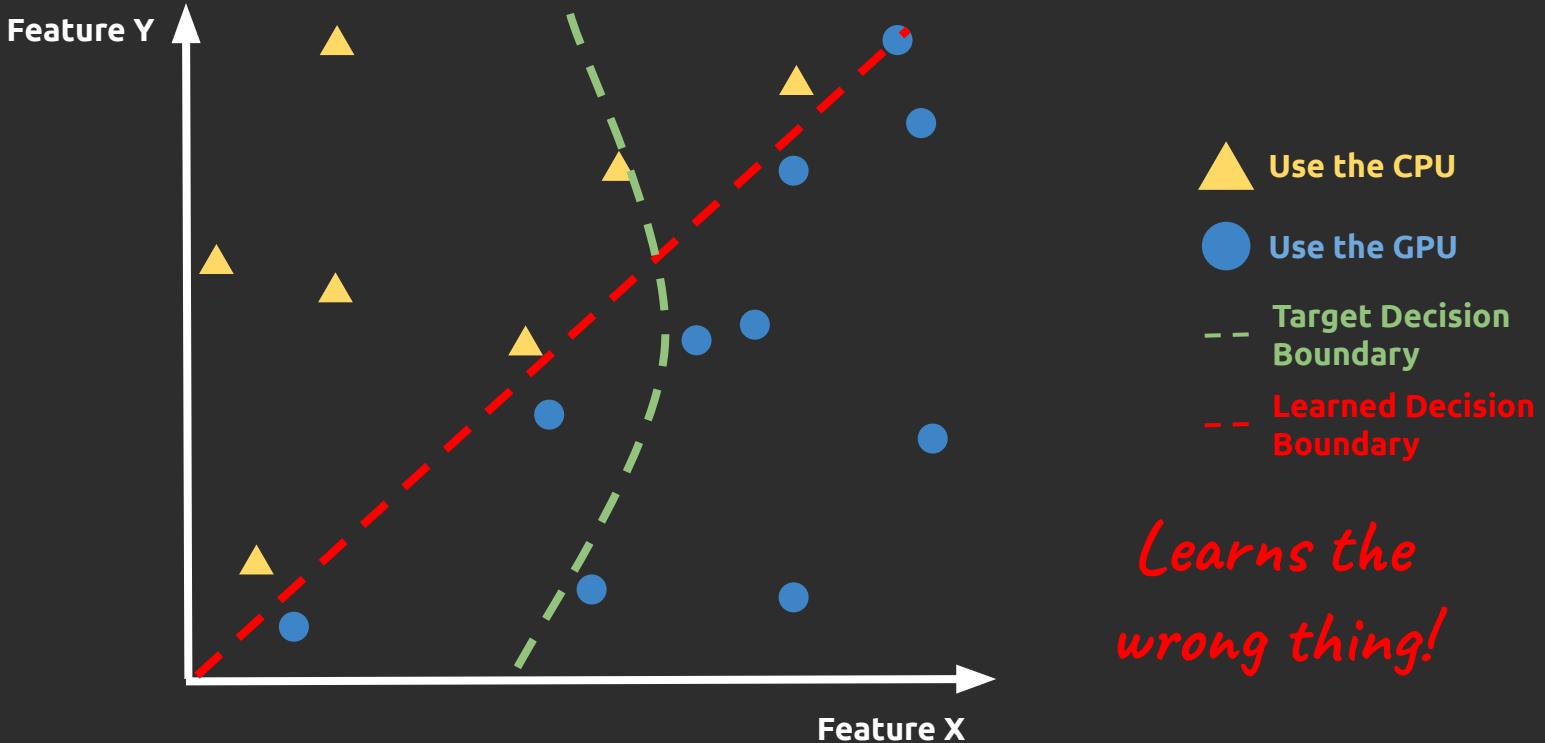
What we have



avg. compiler paper uses 17 benchmarks

source: 25 GPGPU papers from 2013-2016 in CGO, PACT, PPopp HiPC

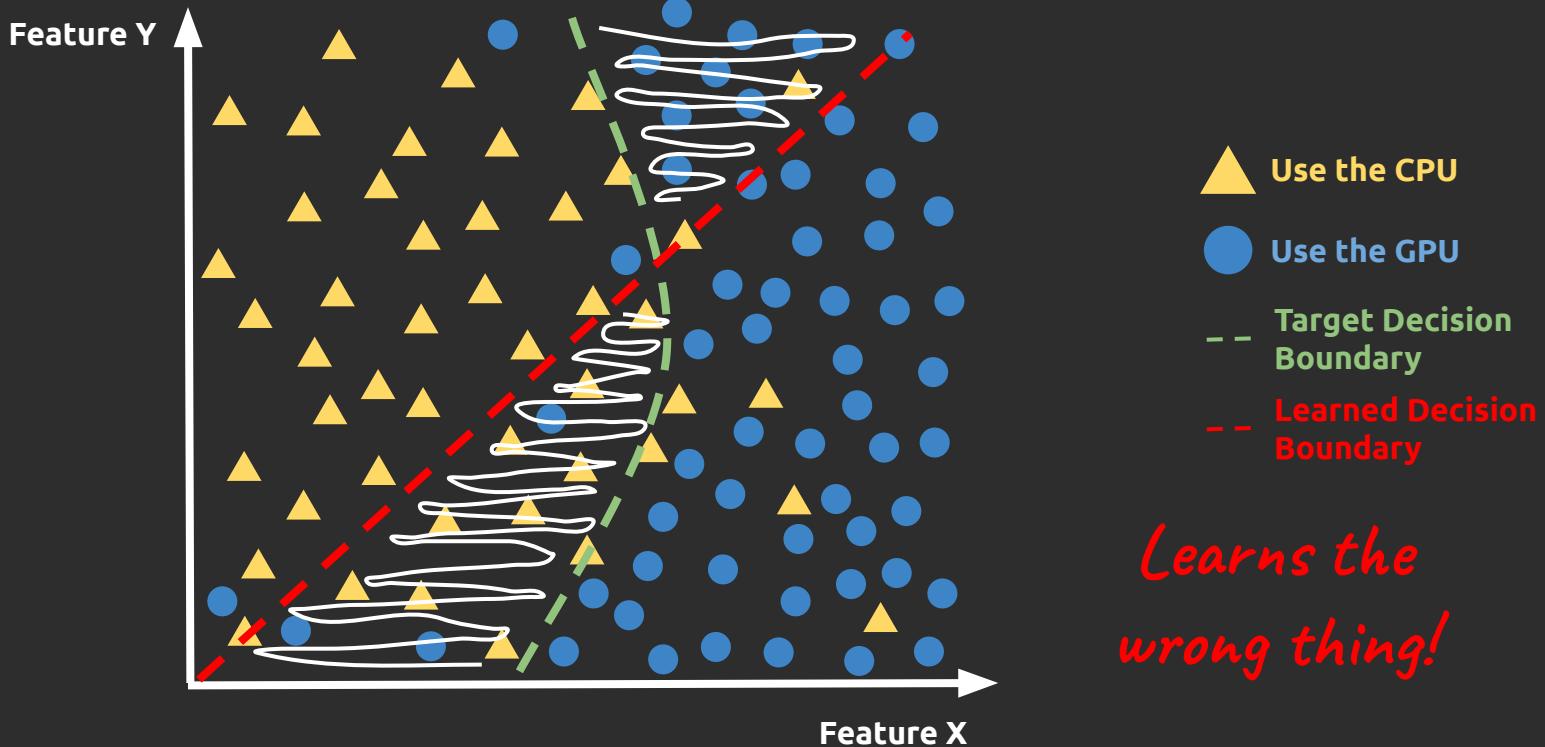
What we have



avg. compiler paper uses 17 benchmarks

source: 25 GPGPU papers from 2013-2016 in CGO, PACT, PPopp HiPC

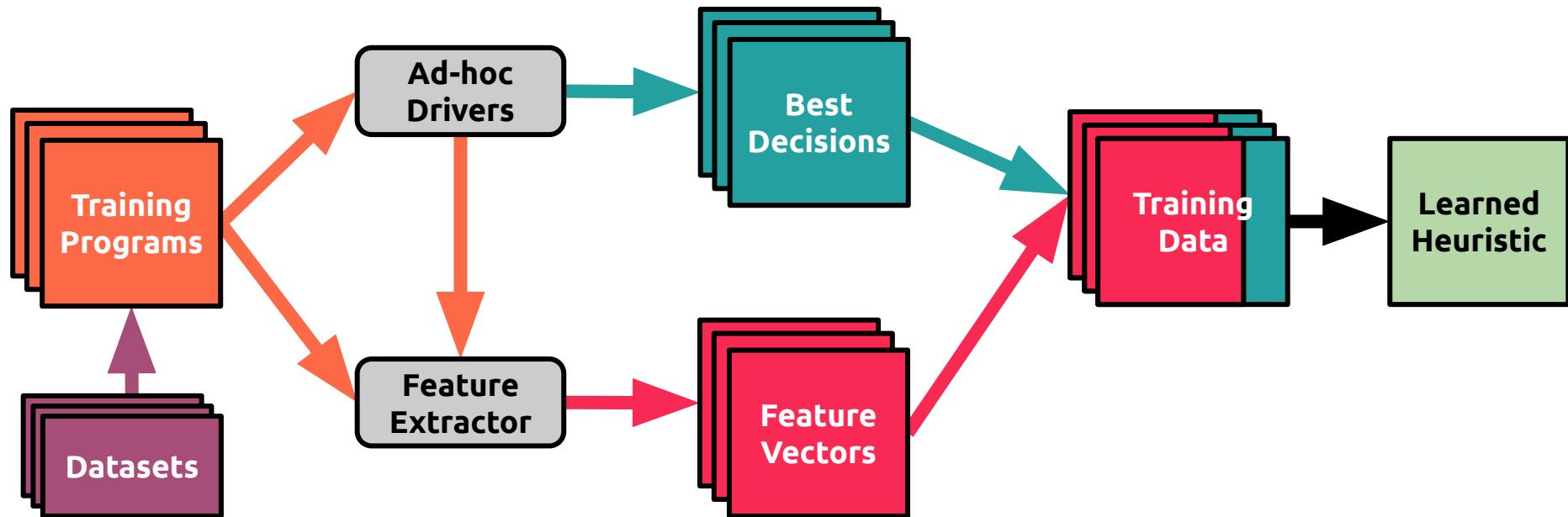
What we have



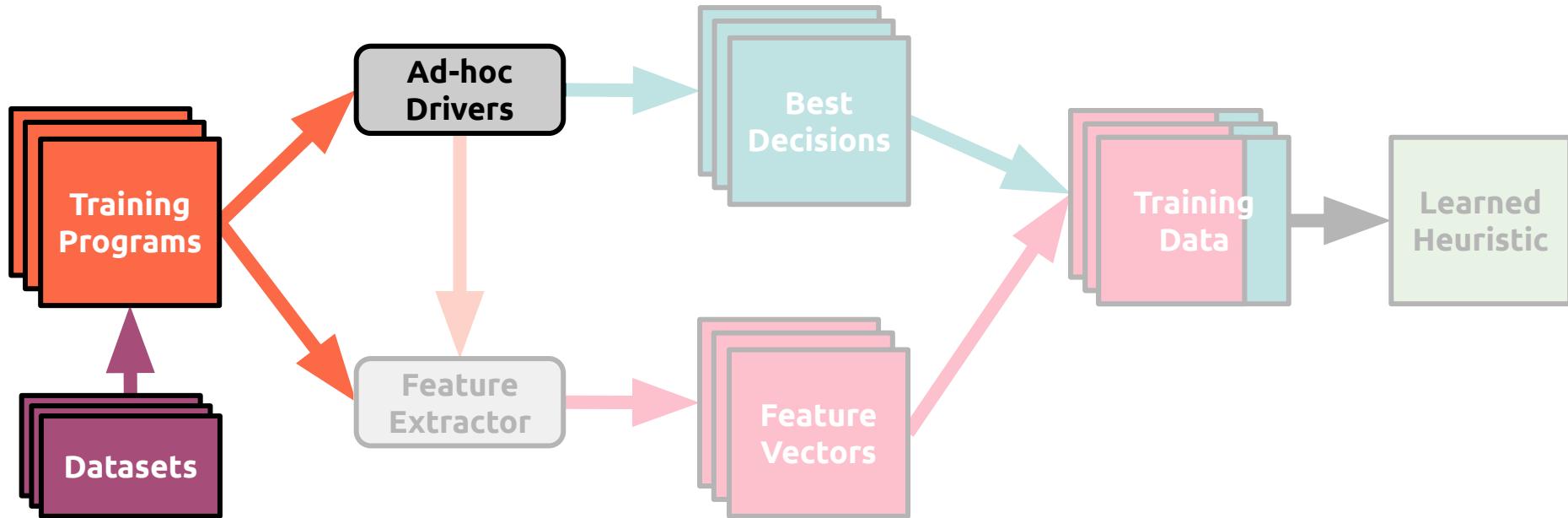
avg. compiler paper uses 17 benchmarks

source: 25 GPGPU papers from 2013-2016 in CGO, PACT, PPopp HiPC

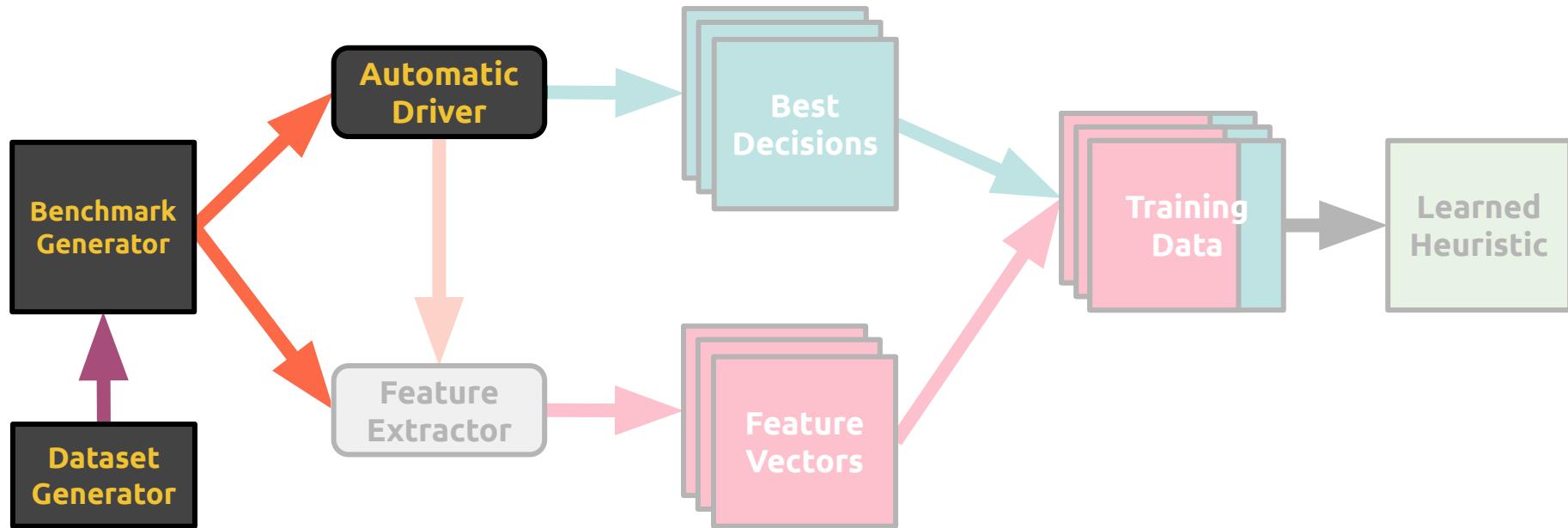
Machine learning for compilers



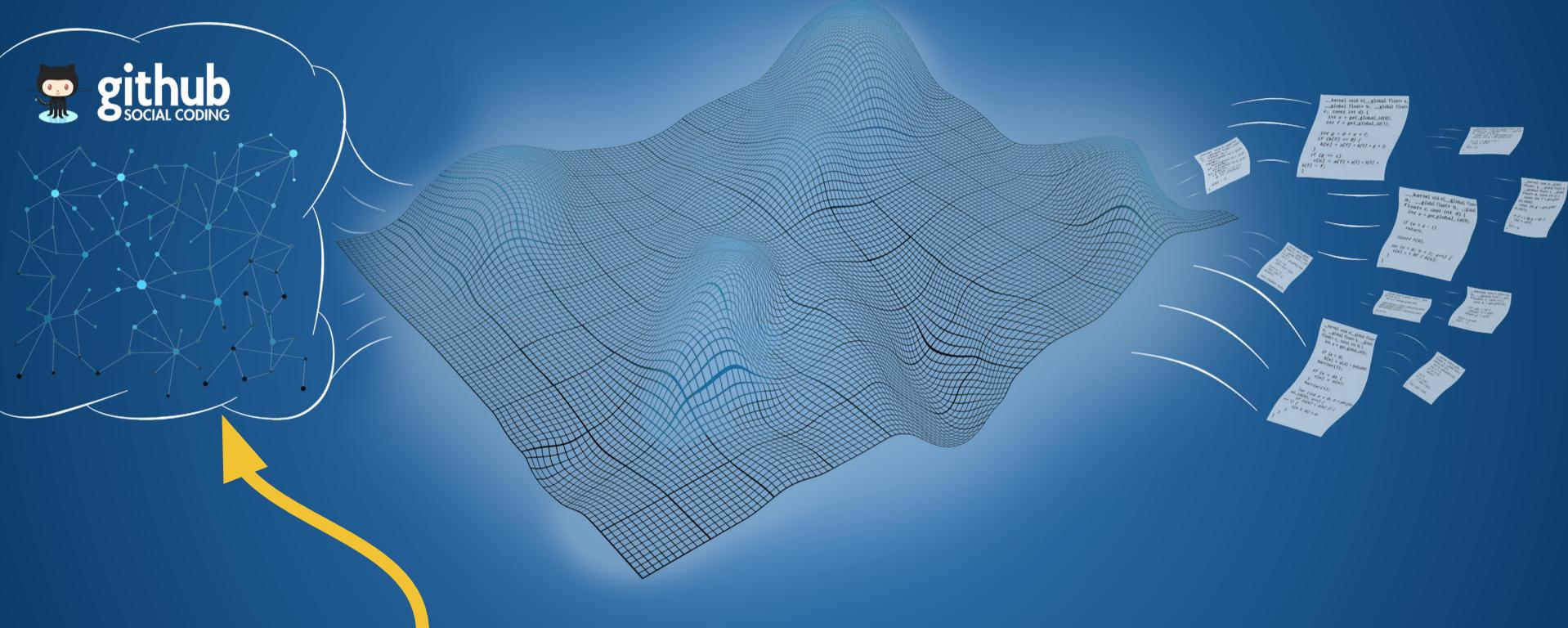
What we have



What we need

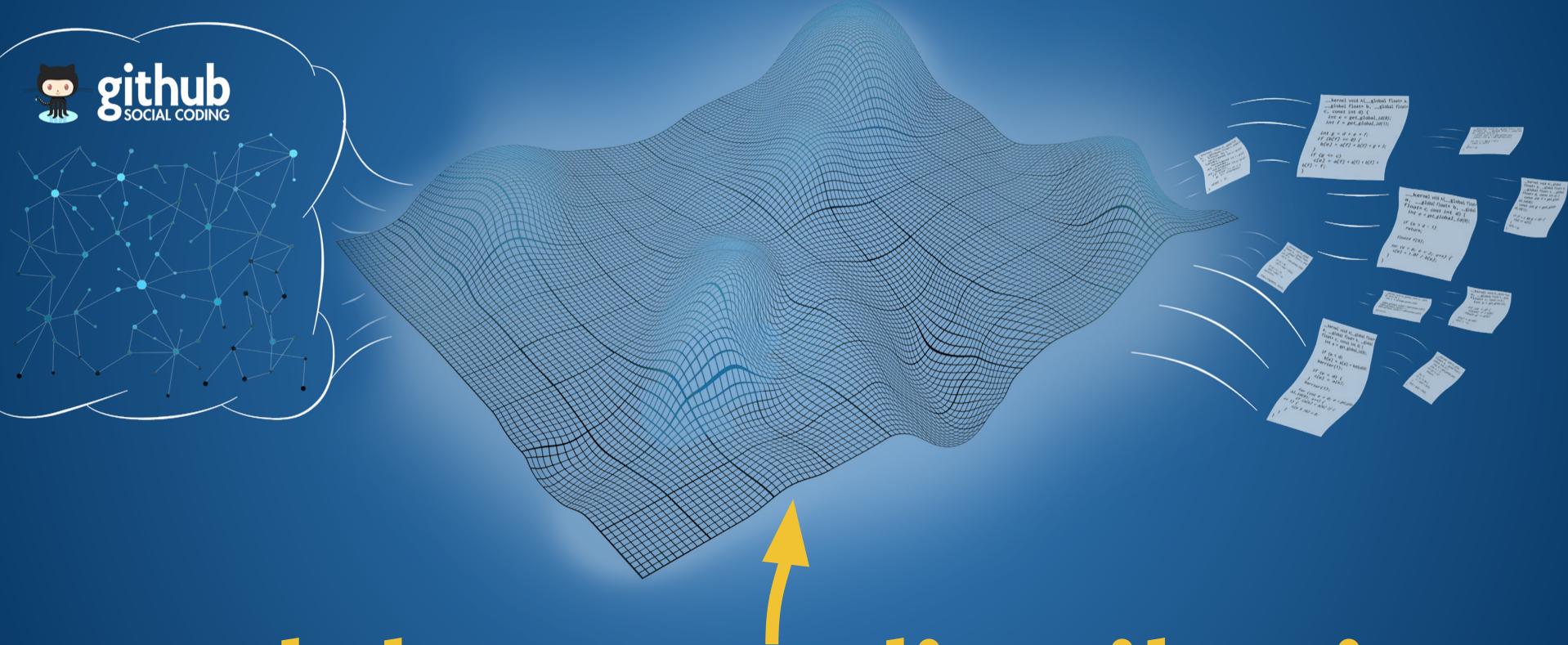


Generating compiler benchmarks



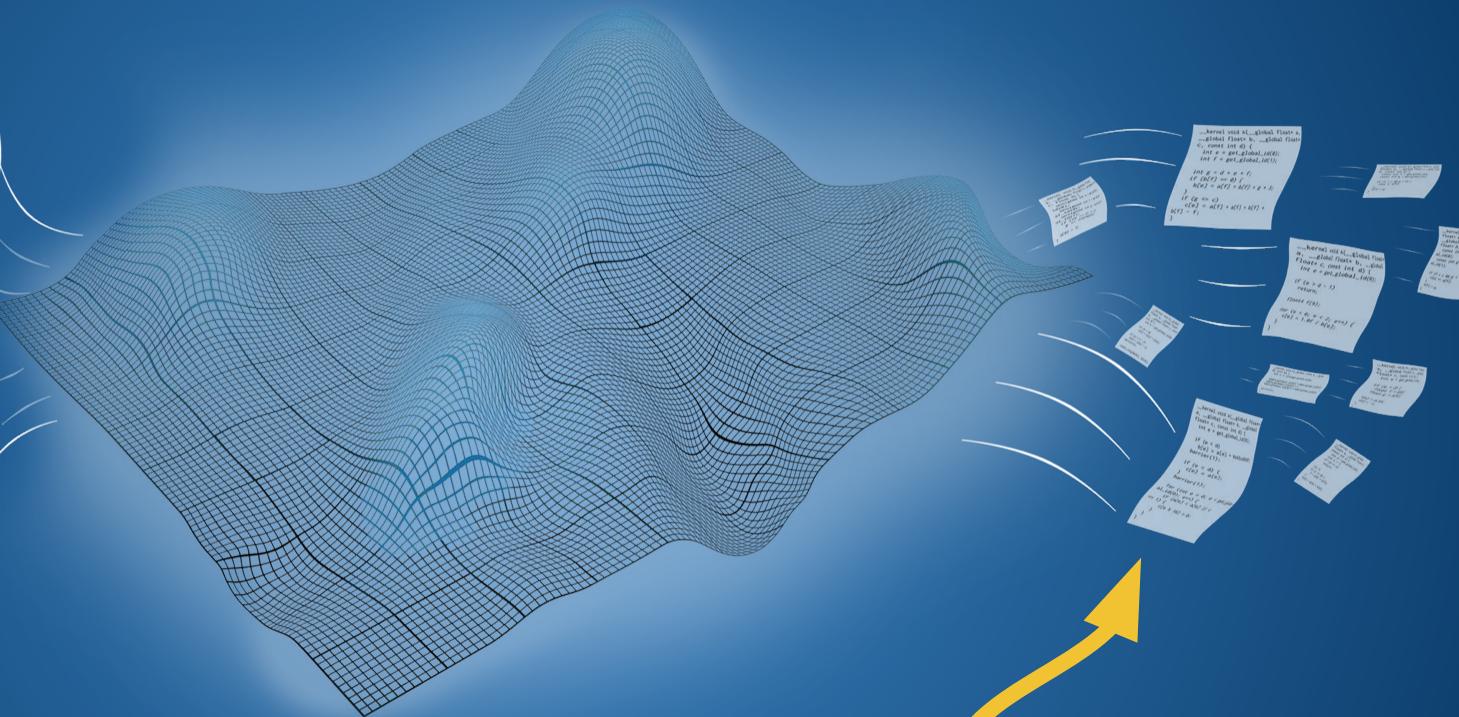
mine code from the web

Generating compiler benchmarks

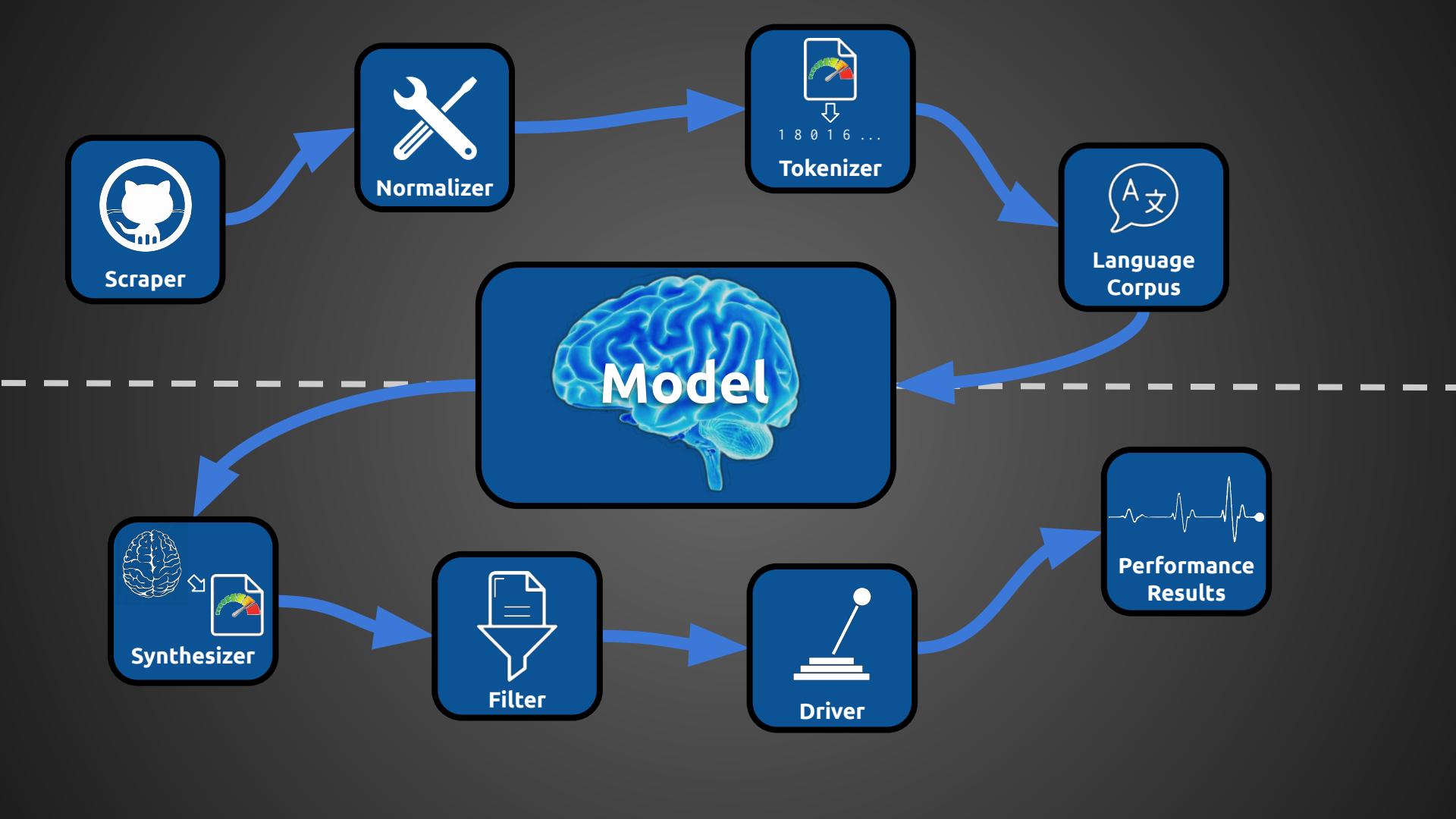


model source distribution

Generating compiler benchmarks



sample language model



Listing 3: Sample 3

```
1 __kernel void A(__global int* a, __global int* b, __global int* c,
2   ↪ __global int* d, const uint e) {
3   const uint f = get_global_id(0);
4
5   if (e == 0 && f == 0)
6     *d = 0;
7   else if (f < e) {
8     int g = b[f];
9     uint h = c[f];
10    if (g > 0) {
11      a[h] = f;
12      h++;
13    }
14    if (f == e - 1)
15      *d = h;
16 }
```

Listing 4: Sample 4

```
1 __kernel void A(__global float* a, __global float* b, __global float* c,
2   ↪ const int d) {
3   int e = get_global_id(0);
4
5   if (e < d) {
6     float f = b[e];
7     float g = a[e];
8     a[e] = f * 3.141592f / (f + 1.0f + e * 1024 - f) - (0.55f * g * 1.0f
9       ↪ / 18.0f + e / 2.0f);
10
11   for (e = 0; e < 30; e++) {
12     c[e] = 0;
13   }
```

Listing 4: Sample 4

```
1 __kernel void A(int a, __global float* b, __global int* c, __global int*
2   ↪ d, __local int* e, int f) {
3   int g = get_local_id(0);
4   e[get_local_id(0)] = 0;
5   barrier(1);
6   while (g < f) {
7     int h = c[g];
8
9     if (h != -1) {
10       __global float* i = b + g * a;
11       float j = 0;
12       int k = 0;
13
14       for (int l = 0; l < a; ++l) {
15         float m = i[l];
16         if (m < 0) {
```

52%
blind test

Listing 7: Sample 7

```
1 __kernel void A(int a, int b, int c, __global const float* d, __global
2   ↪ const float* e, __global float* f, float g) {
3   const int h = get_local_id(0);
4   const int i = get_group_id(0);
5
6   const int j = 4 * i + h;
7   const int k = 4 * i + h + a;
8   if (4 * i + h + a < c) {
9     float l = 0.0;
10    float m = 0.0;
11    float n = 0.0;
12
13    const float o = d[3 * (4 * i + h + a)];
14    const float p = d[3 * (4 * i + h + a) + 1];
15    const float q = d[3 * (4 * i + h + a) + 2];
16
17    for (int r = 0; r < c; r++) {
18      const float s = d[3 * r] - o;
19      const float t = d[3 * r + 1] - p;
20      const float u = d[3 * r + 2] - q;
21      const float v = (d[3 * r] - o) * (d[3 * r] - o) + (d[3 * r + 1] -
22        ↪ p) * (d[3 * r + 1] - p) + (d[3 * r + 2] - q) * (d[3 * r +
23        ↪ 2] - q) + g;
24
25      const float w = e[r] / (((d[3 * r] - o) * (d[3 * r] - o) + (d[3 *
26        ↪ r + 1] - p) * (d[3 * r + 1] - p) + (d[3 * r + 2] - q) * (d
27        ↪ [3 * r + 2] - q) + g) * sqrt((d[3 * r] - o) * (d[3 * r] - o)
28        ↪ + (d[3 * r + 1] - p) * (d[3 * r + 1] - p) + (d[3 * r + 2] -
29        ↪ q) * (d[3 * r + 2] - q) + g));
30
31      m = m + (d[3 * r + 1] - p) * w;
32      n = n + (d[3 * r + 2] - q) * w;
33    }
34
35    f[3 * (4 * i + h)] = l;
36    f[3 * (4 * i + h) + 1] = m;
37    f[3 * (4 * i + h) + 2] = n;
38  }
39 }
```

Listing 10: Sample 10

```
1 __kernel void A(__global ulong *a) {
2   int i, j;
3   struct S0 c_8;
4   struct S0* p_7 = &c_8;
5   struct S0 c_9 = {
6     {0x43250E6DL, 2UL}, {0x43250E6DL, 2UL}, {0x43250E6DL, 2UL},
7     {0x43250E6DL, 2UL}, {0x43250E6DL, 2UL}, {0x43250E6DL, 2UL},
8     {0x43250E6DL, 2UL}, {0x43250E6DL, 2UL}, {0x4BF90EDCAD2086BDL,
9
10    };
11   c_8 = c_9;
12   barrier(0 | 1);
13   B(d);
14   barrier(0 | 1);
15   uint64_t e = 0xFFFFFFFFFFFFFFFUL;
16   a[get_linear_global_id()] = e ^ 0xFFFFFFFFFFFFFFFUL;
17 }
```

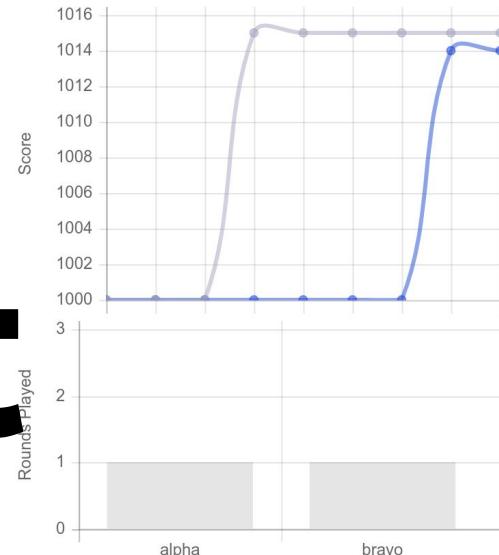
96% control group

Round 1

Player: 973, Robot: 1028

```
_kernel void A(__global __write_only int *a, __global __read_only int *b,
               const float c, const float d, const float e, const float f,
               const int g, const int h, const int i) {
    int j = get_global_id(0);
    int k = get_global_id(1);
    float l = c + j * e / g;
    float m = d + k * f / h;
    int n = 0;
    float o = 0;
    float p = 0;
    while (n < i) {
        float q = o * o;
        float r = p * p;
        if ((q + r) > 4.0) break;
        float s = q - r + l;
        p = 2 * o * p + m;
        o = s;
        ++n;
    }
    if (n < i)
        a[k * g + j] = b[n];
    else
        a[k * g + j] = 0xff000000;
}
```

Try it



Human

Robot

Synthesizing Benchmarks for Predictive Modeling

C. Cummins, P. Petoumenos, Z. Wang, H. Leather

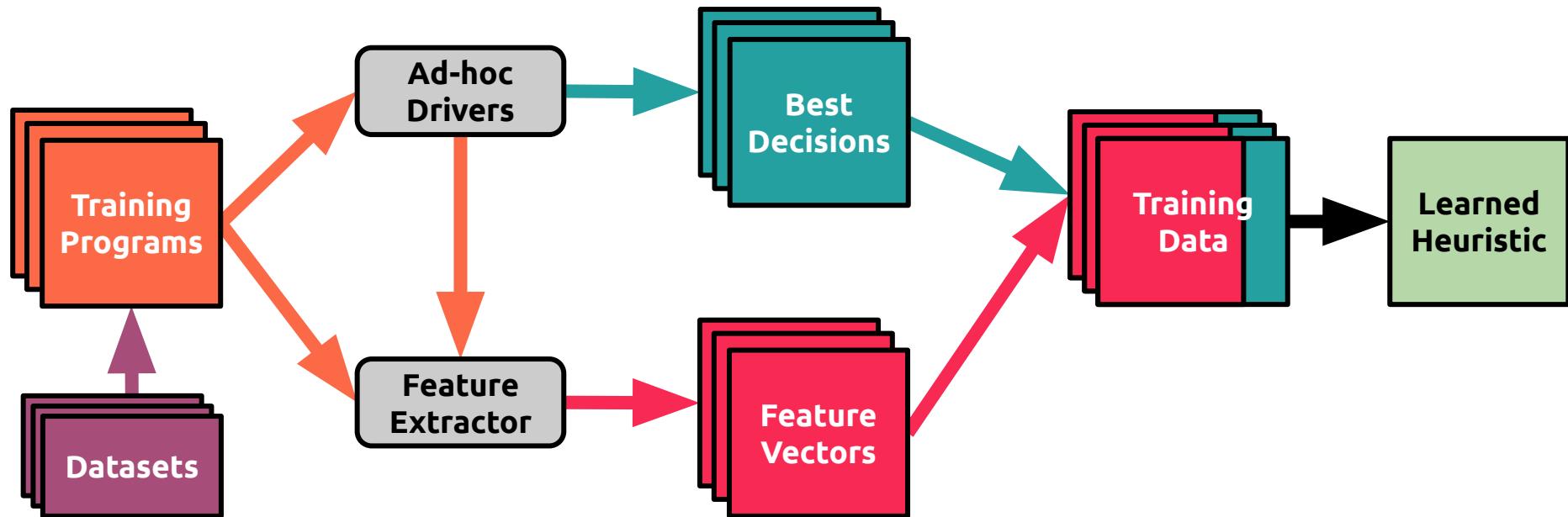


DL approach to model PL **syntax & usage**

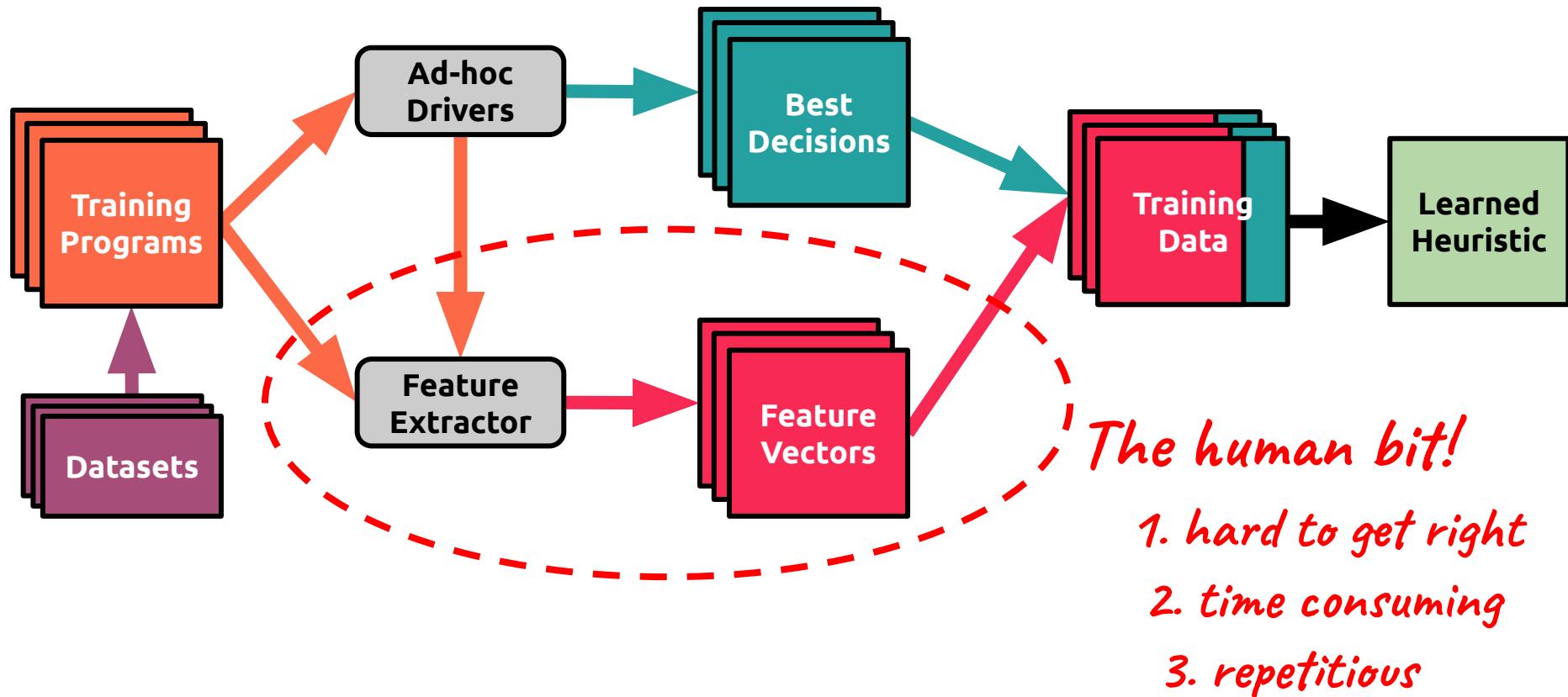
First **general purpose** benchmark generator

4.30x better model **performance**

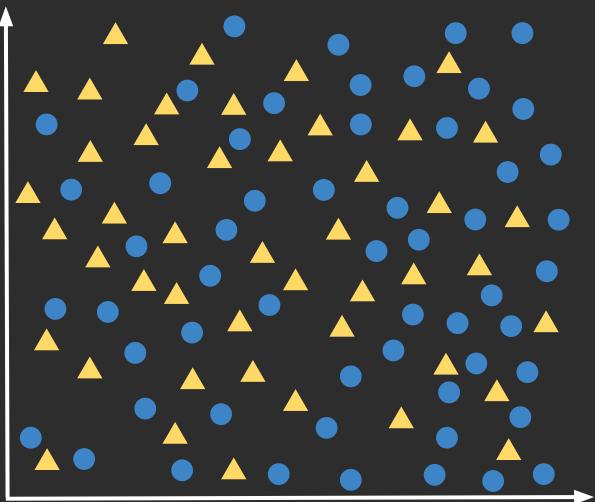
Machine learning for compilers



Machine learning for compilers

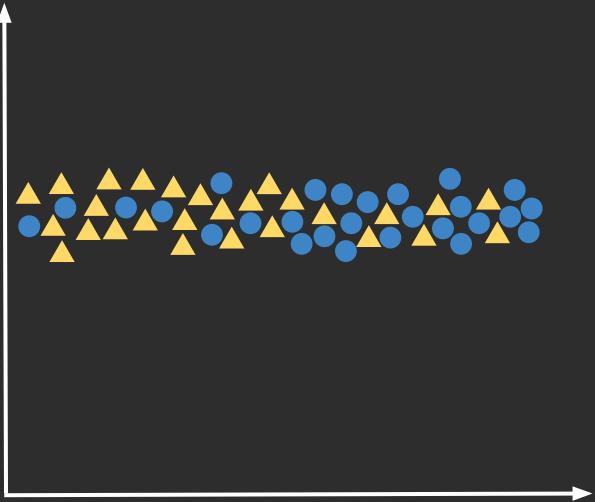


Ways to fail



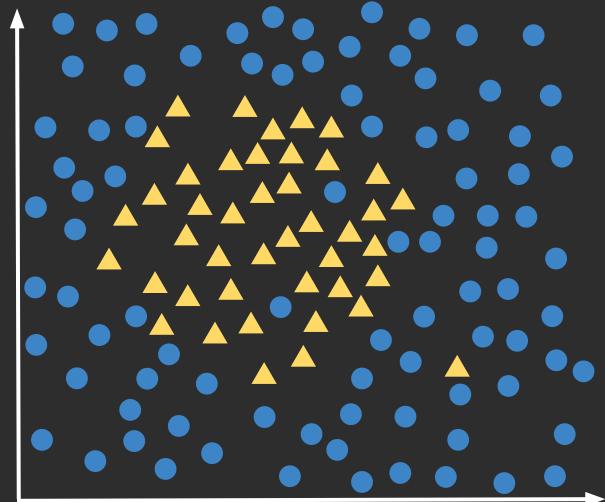
Irrelevant

e.g. not capturing the right information



Incomplete

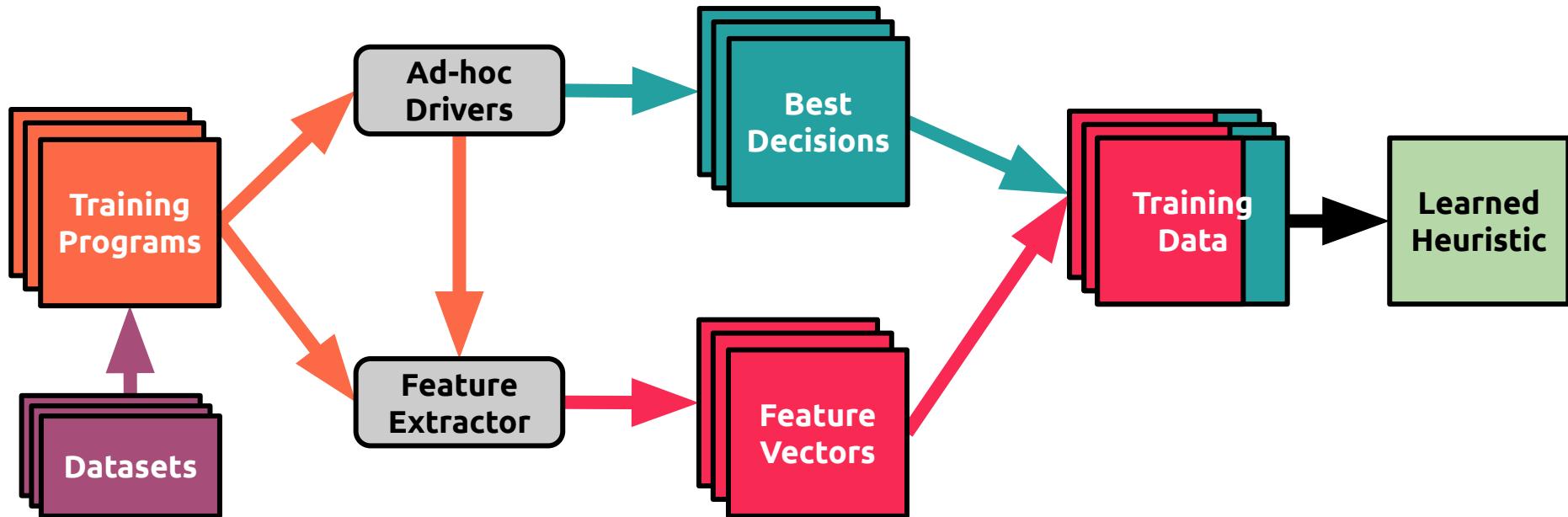
e.g. missing critical information



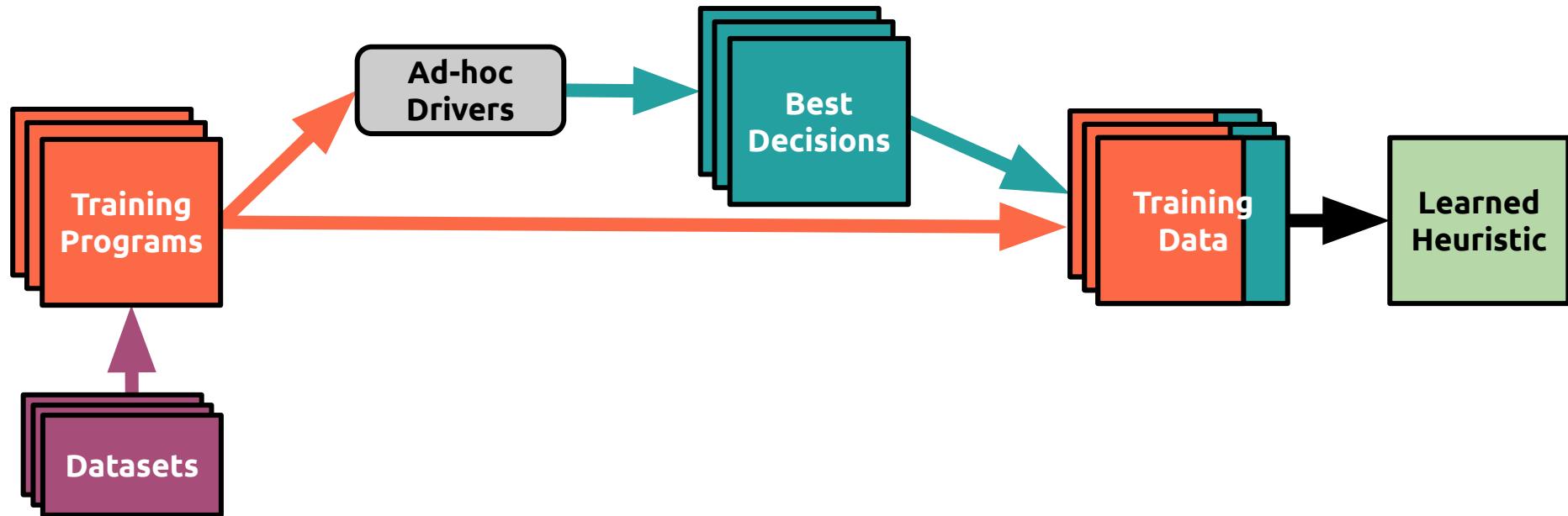
Unsuitable

e.g. wrong combination of features+model

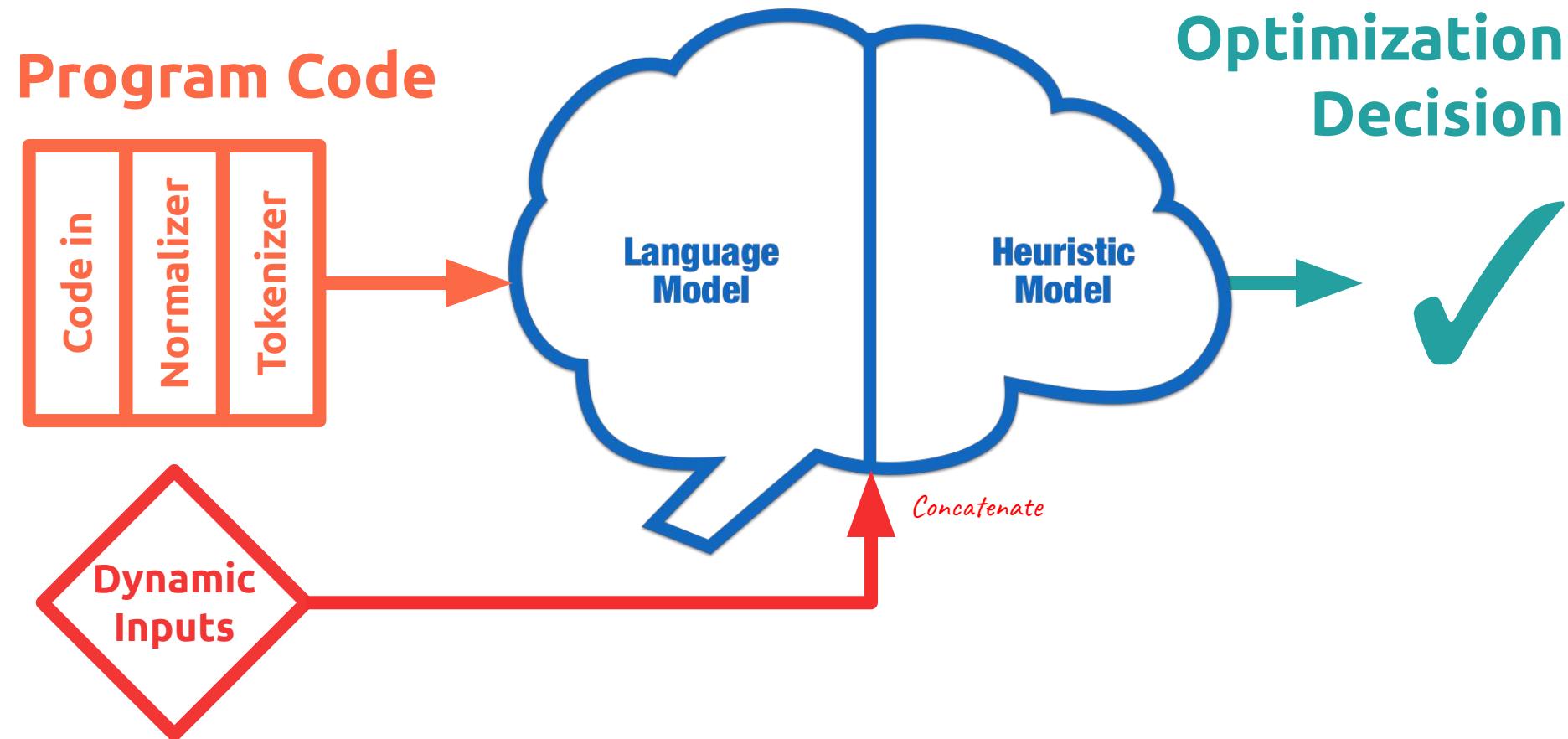
What we have



What we need

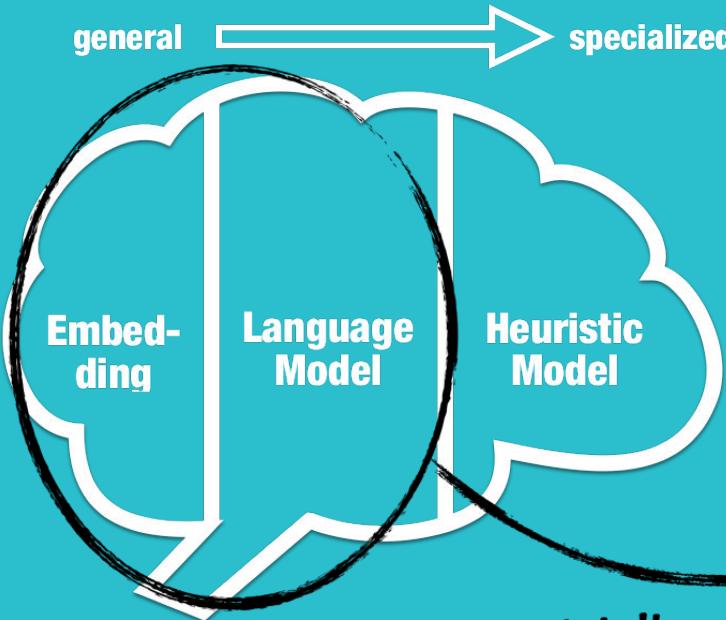


End-to-end optimization learning

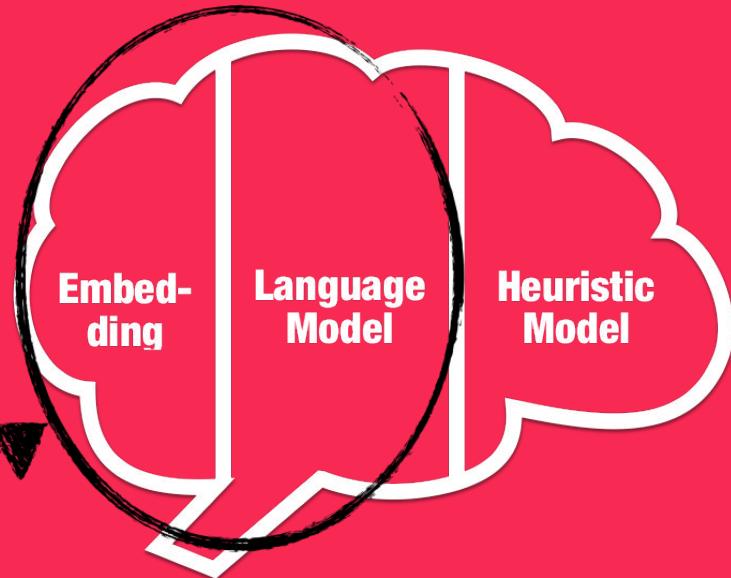


Learning across optimizations

Optimization A



Optimization B



initialize with values
+ fine-tune

End-to-end Deep Learning of Optimization Heuristics

C. Cummins, P. Petoumenos, Z. Wang, H. Leather



Heuristics **without** features

12.5% **better** than using features

Learning **across** heuristics

Compilers break



Compile-time failure

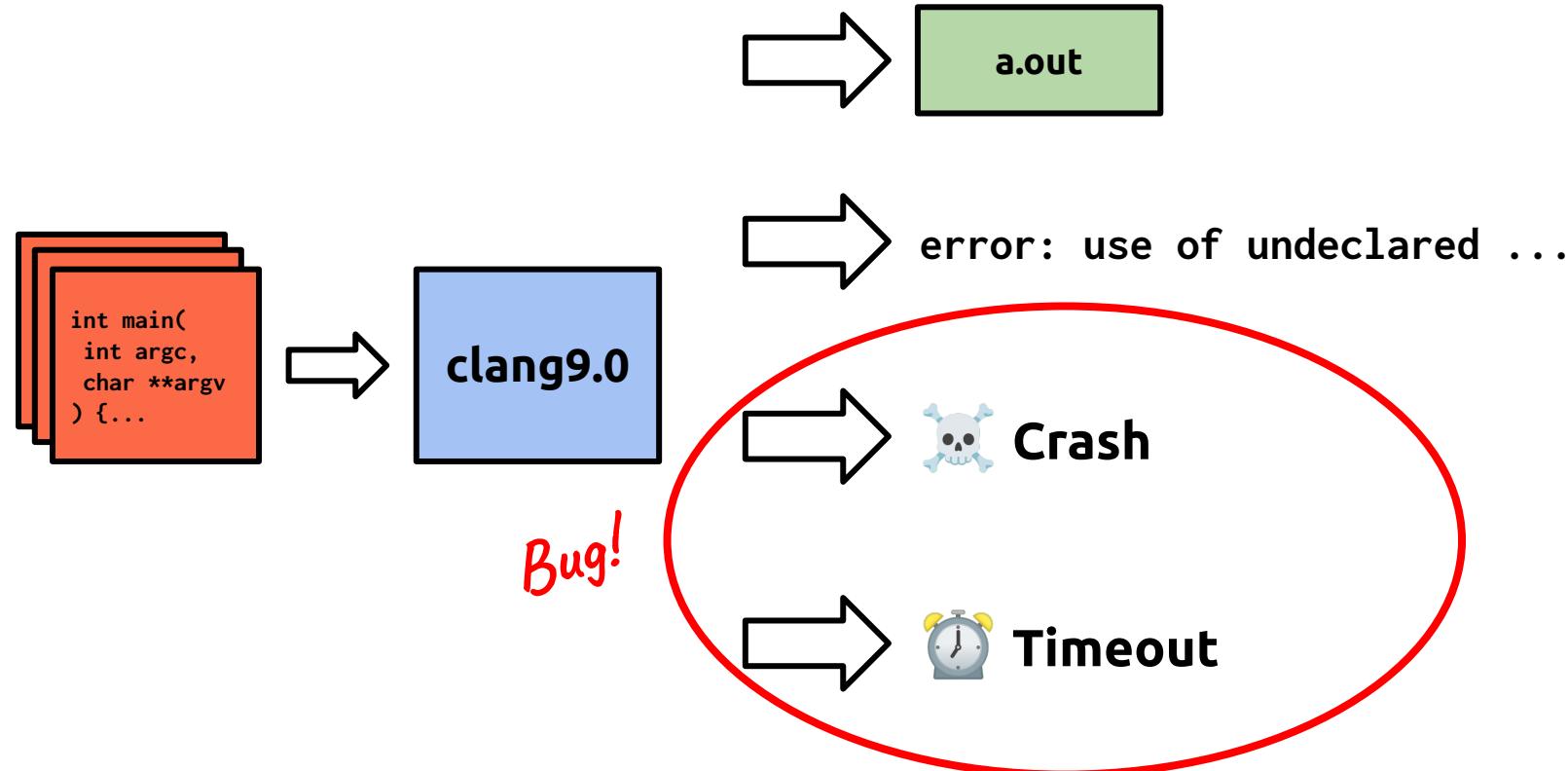


Broken software

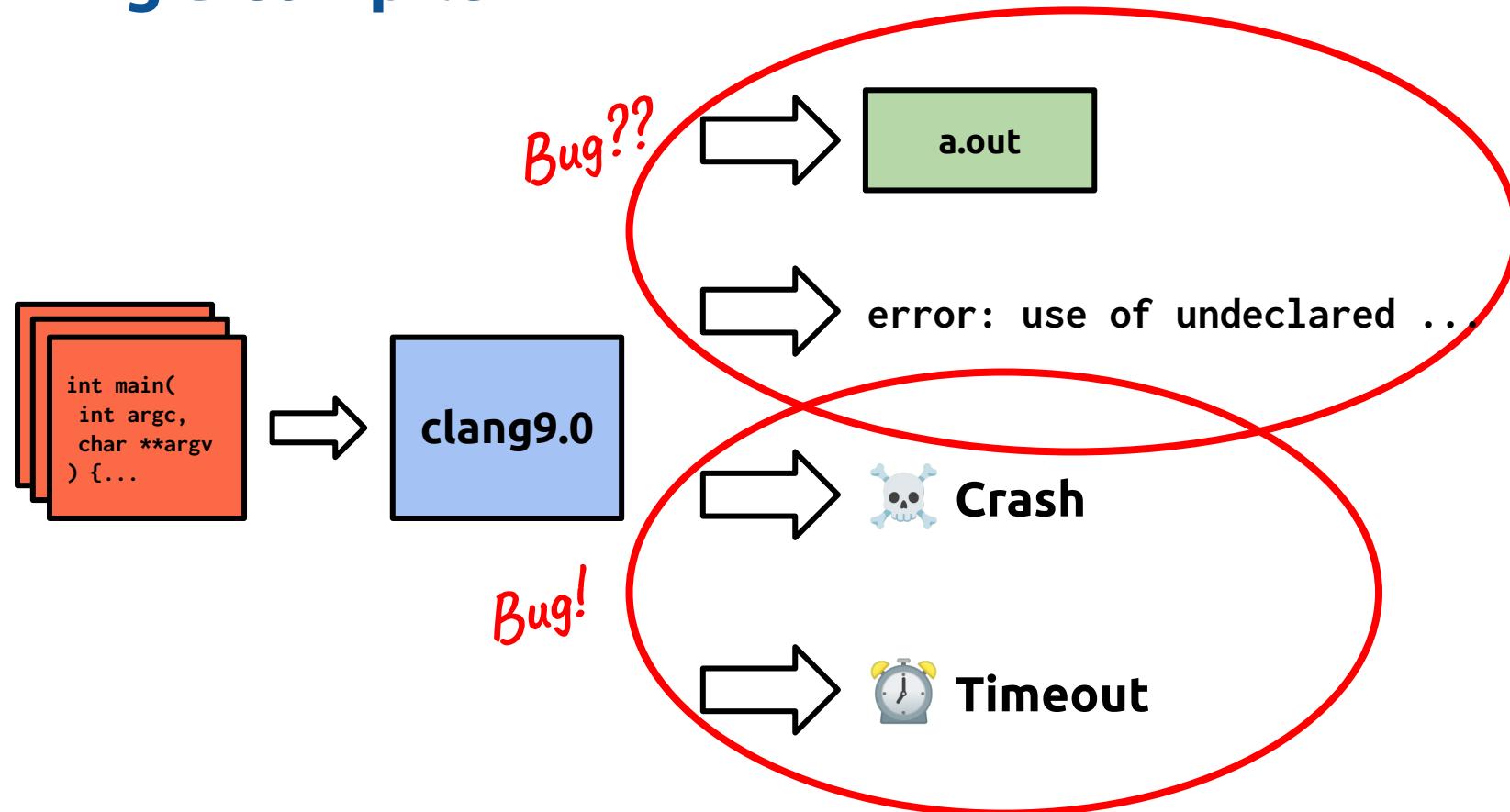
Regression suites

1. Slow
2. Late
3. Expensive
4. Incomplete

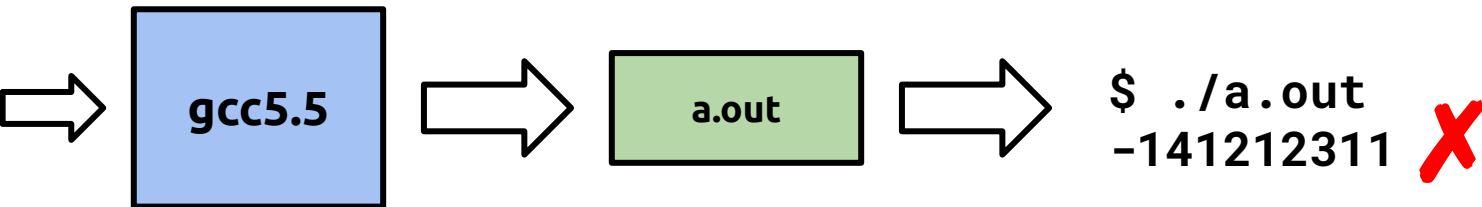
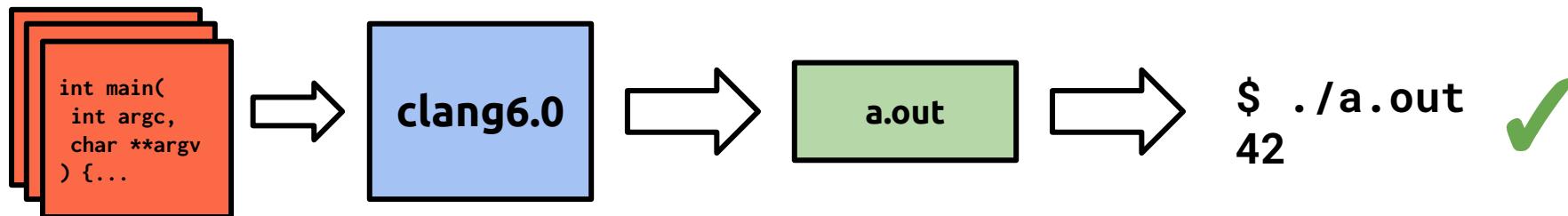
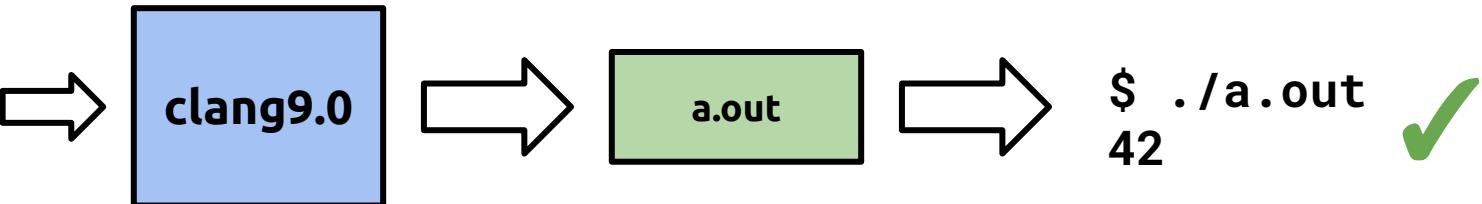
Fuzzing a compiler



Fuzzing a compiler

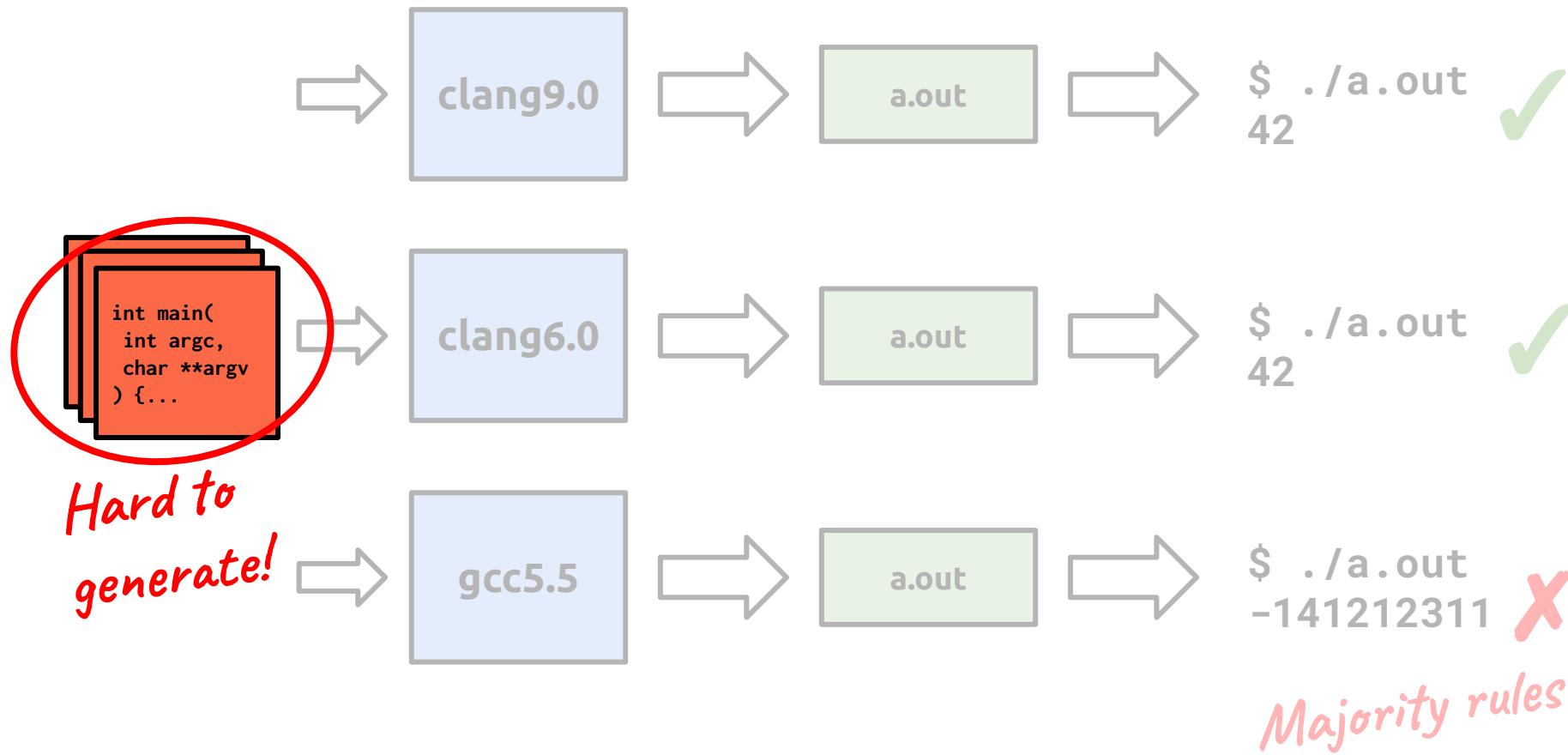


Differential testing compilers

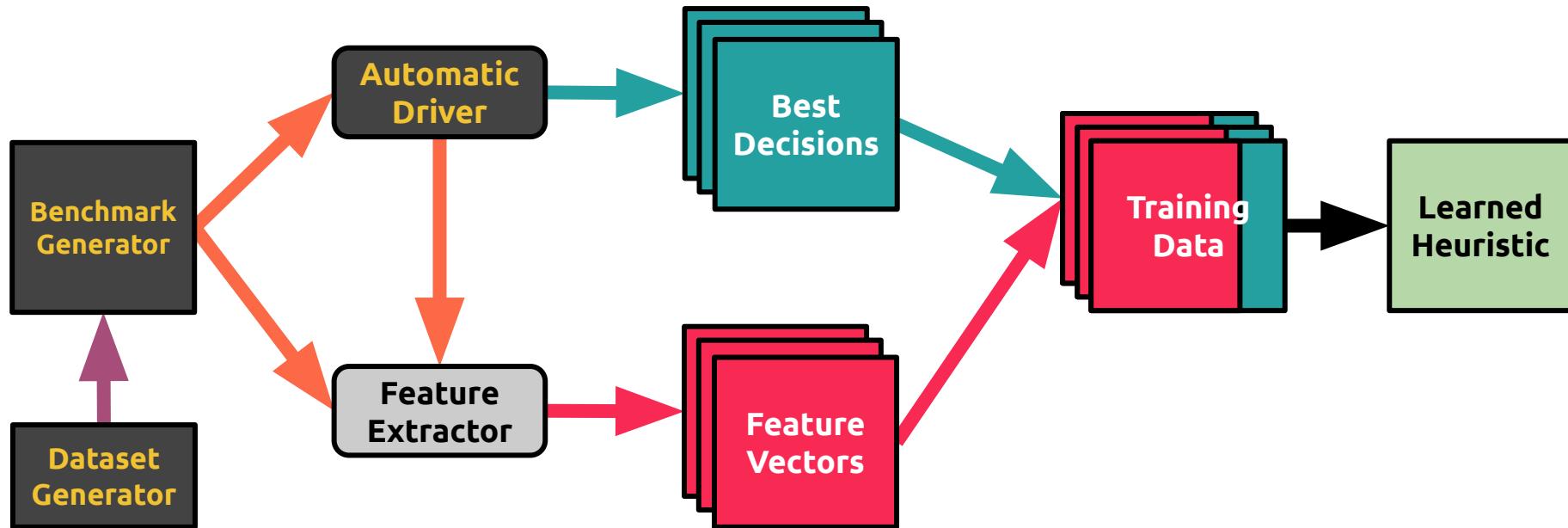


Majority rules

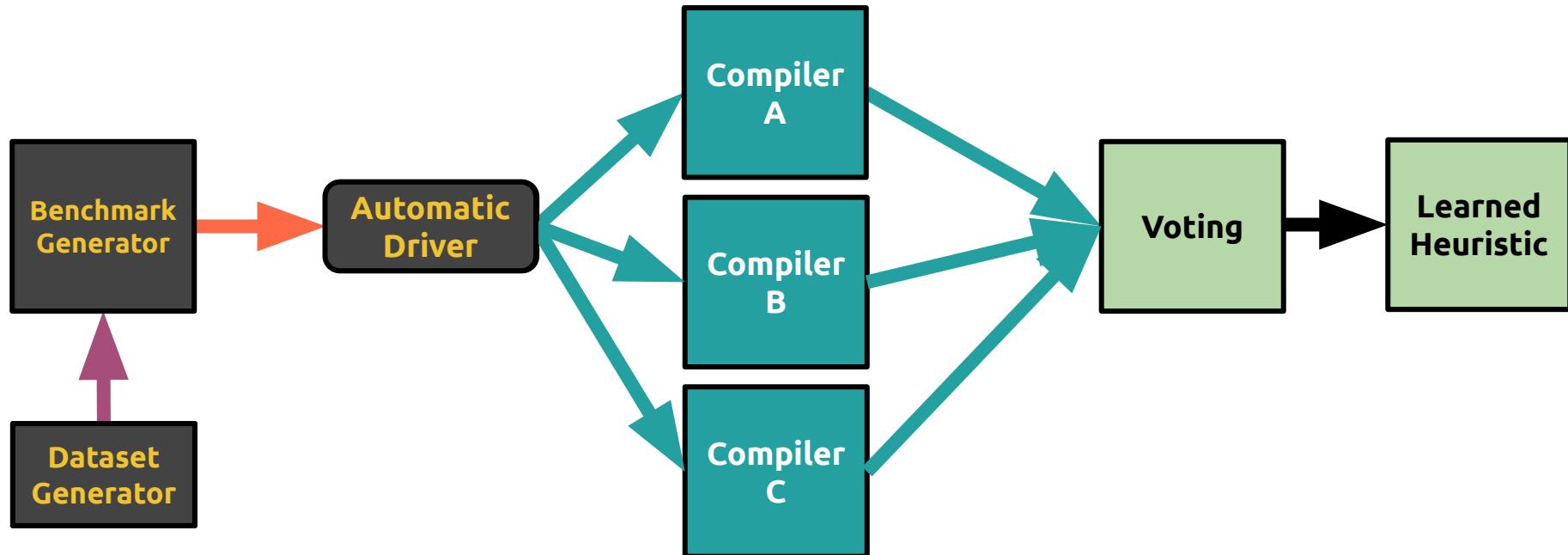
Differential testing compilers



What we have



What we need



Compiler Fuzzing through Deep Learning

C. Cummins, P. Petoumenos, Z. Wang, H. Leather



Automatic inference of fuzzers

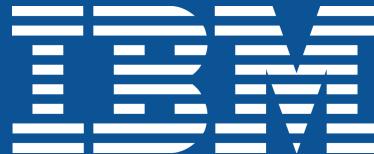
102× less code than state-of-art

Small, interpretable test cases

In production



In prototype



Thesis

1.  **Can we automate benchmark engineering?**
2.  **Can we automate feature engineering?**
3.  **Can we use (1) to automate testing?**

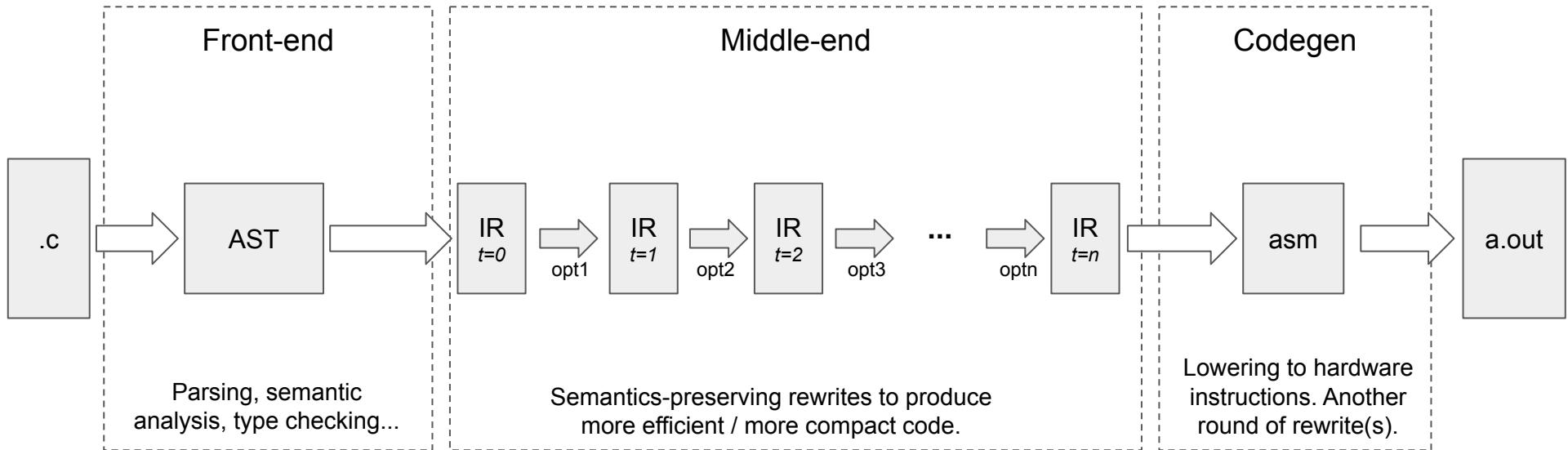
overview

- 1. What I used to do**
- 2. What I now do**
- 3. What you should do**

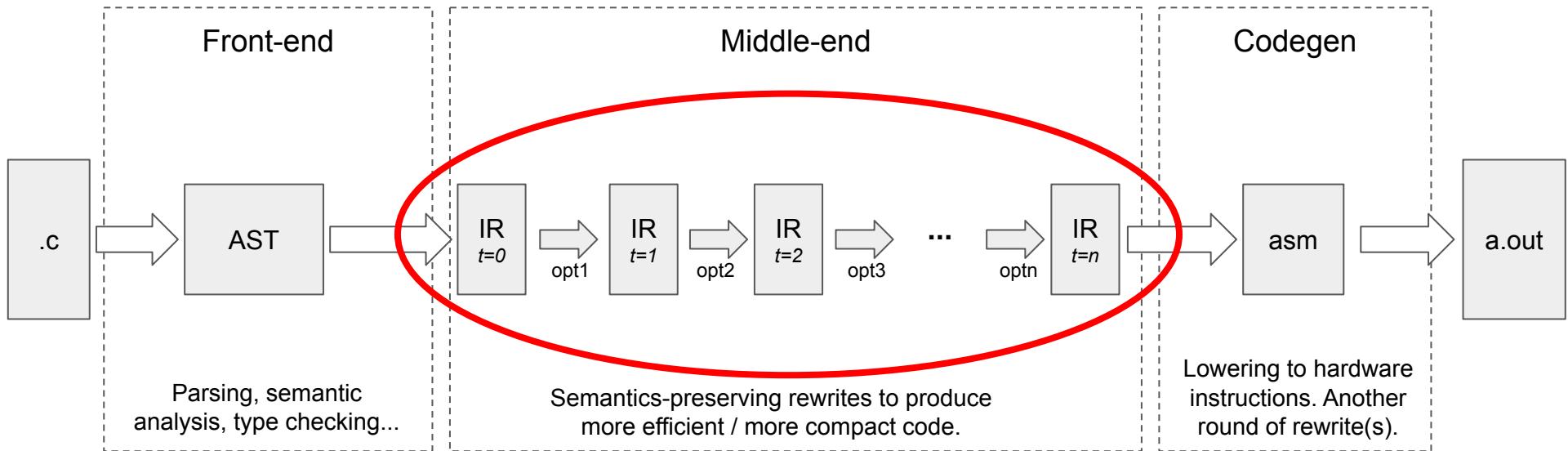
overview

1. What I used to do
2. What I now do
3. What you should do

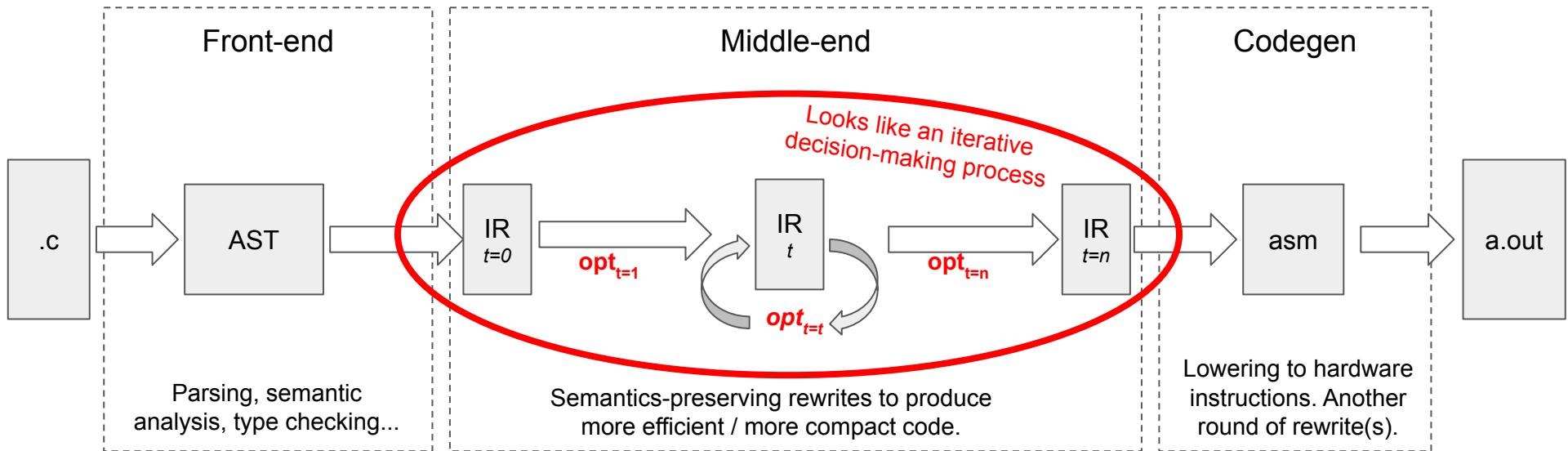
High-level overview of compiler



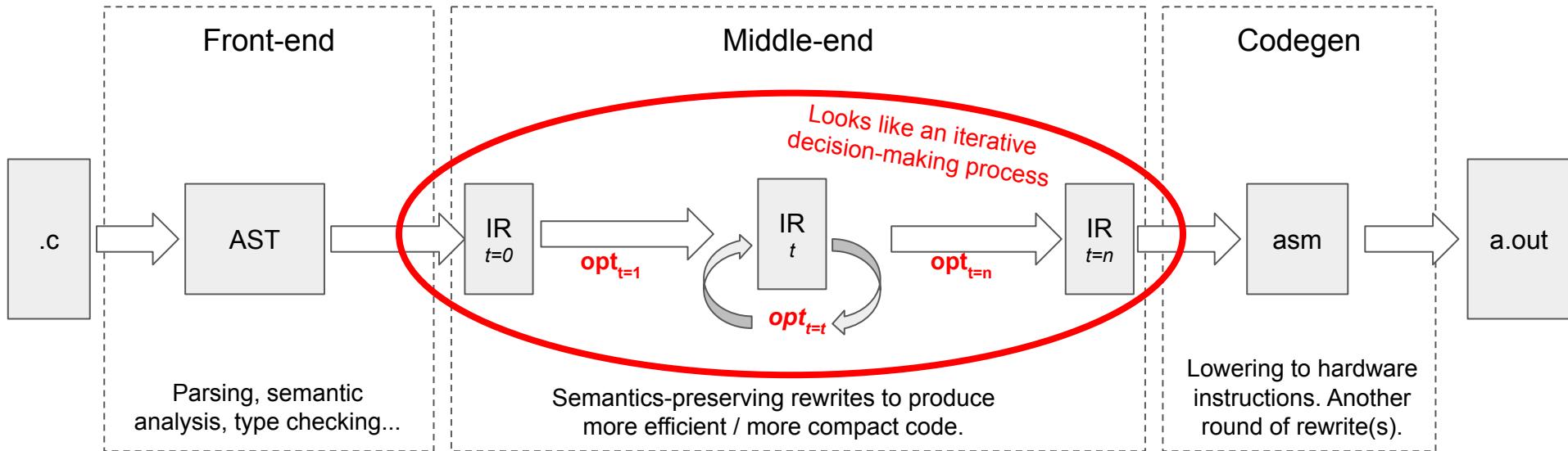
Optimizations are applied in a fixed order



Choosing the best is known as "Phase Ordering"



Choosing the best is known as "Phase Ordering"



Search shows large performance wins.
([up to 2.23x](#))

But search is slow, and must be repeated for each program, input, hardware config, compiler version.



facebook

1 Hacker Way

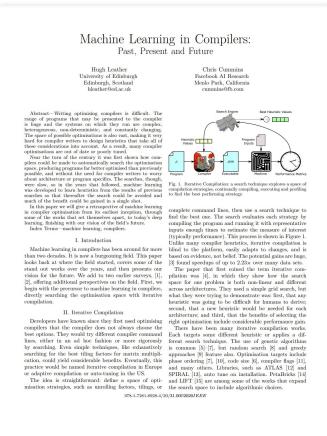
ML for Compilers, recap



Large wins to be had!



Months of work to set up compiler infrastructure + analyses



Short, high-level overview

<https://chriscummins.cc/pub/2020-fd1.pdf>



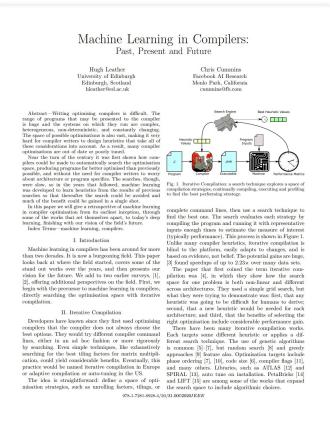
More comprehensive survey

<https://arxiv.org/pdf/1805.03441.pdf>

ML for Compilers, recap

Large wins to be had!

Months of work to set up compiler infrastructure + analyses



Short, high-level overview

Machine Learning in Compiler Optimisation
Zheng Wang and Michael O’Rourke

Zhen Wang and Michael O'Barak

However, in the last decade, machine learning based research has moved from an obscure research niche to a main stream research area. In this article, we describe the relationship between machine learning and research methods and introduce the main concepts of features, models, training and algorithms. We also discuss how machine learning can provide a road map for the wide variety of different research areas. We conclude with a discussion on the use of machine learning in research and its potential to revolutionize the way we approach scientific introduction to the first learning module of machine learning and machine learning as a predictor of the outcome where we find machine learning completion.

Letter to the Editor: Machine Learning, Code Optimize
and L. Introduction

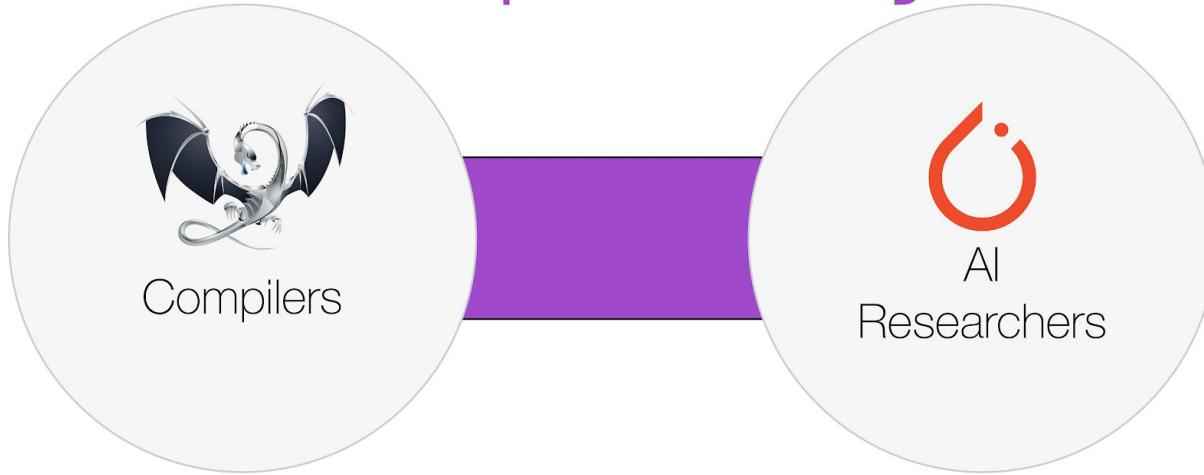
What would you do if you were building a model for a machine learning system? If you are a developer, complex compiler programming languages exist to help you build your system. If you are a scientist or engineer, there are the well-known [13, 11, 17] what-if tools that can help you analyze your system. On the other hand, it is an area of artificial intelligence research to find ways to automatically analyze and optimize machine learning systems. In this letter, we introduce a machine learning system that automatically analyzes machine learning code generated by TensorFlow [1]. All aspects, predicting the size of a certain component, global memory access patterns, and so on, are analyzed by our system. It has as its goal to make the system more efficient. We believe that this kind of analysis will be useful for many applications.

arXiv
research domain.
A. It's all about optimization
Compiler has two jobs – translation and optimization.
Machine learning is part of a tradition in computer science and compilation is increasing automation. The 50s to 70s were spent trying to generate compiler instructions, e.g. for x86.

² Using the same methodology, I was able to observe the effects of the introduction of a new product on the sales of existing products in the same market. The results were published in *Journal of Marketing Research*, 1990, 27, 121-132.

More comprehensive survey

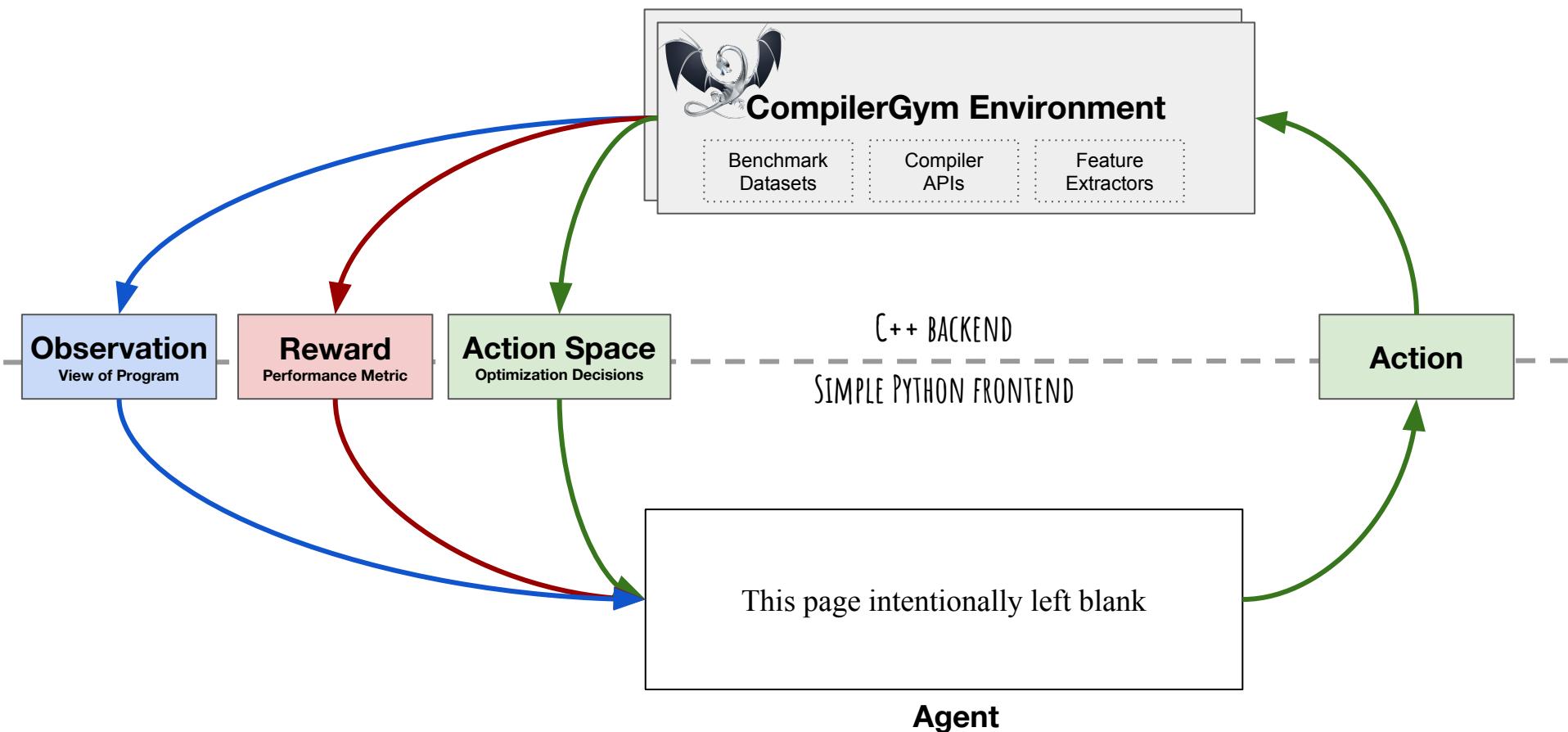
CompilerGym

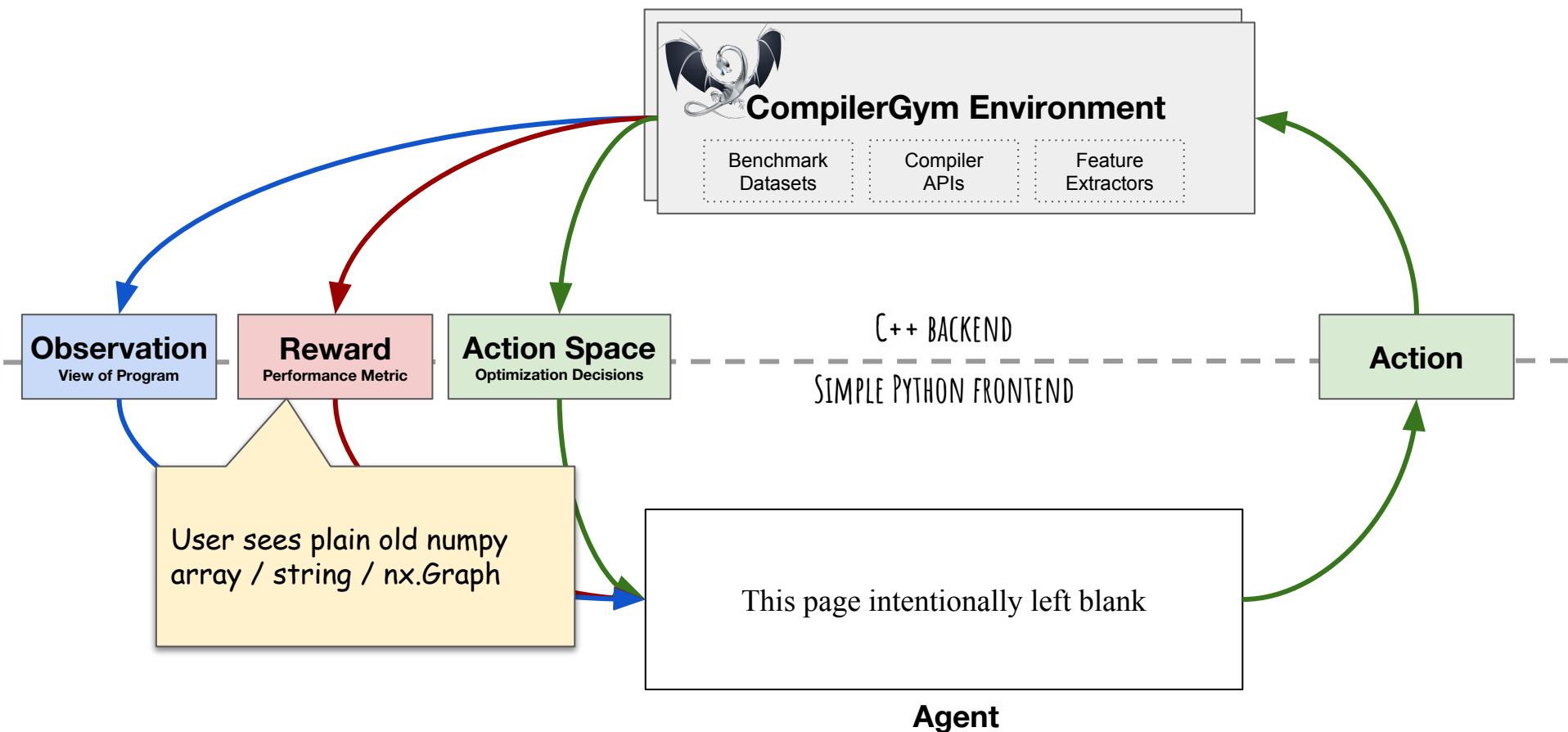


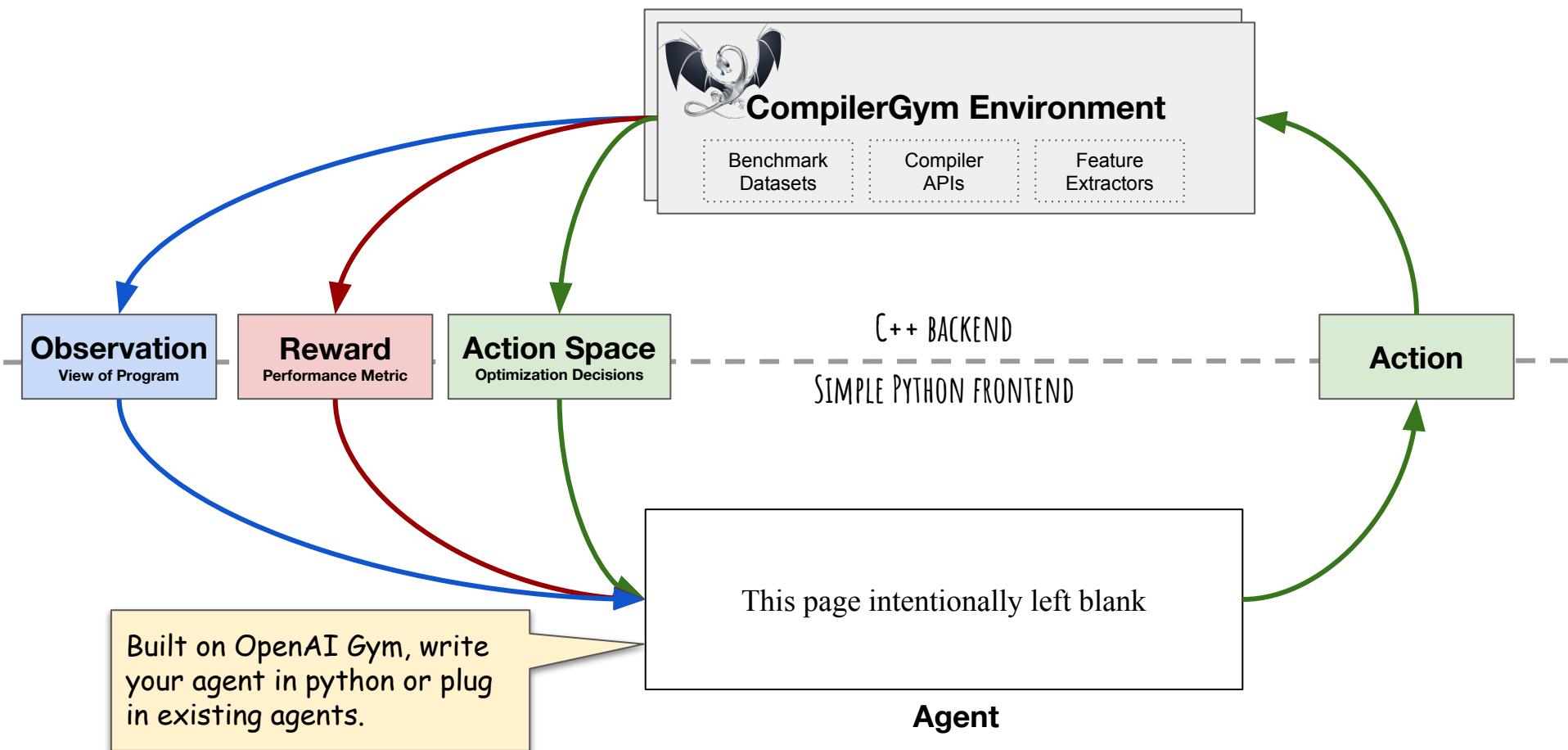
Making Compiler Optimization Accessible to All.

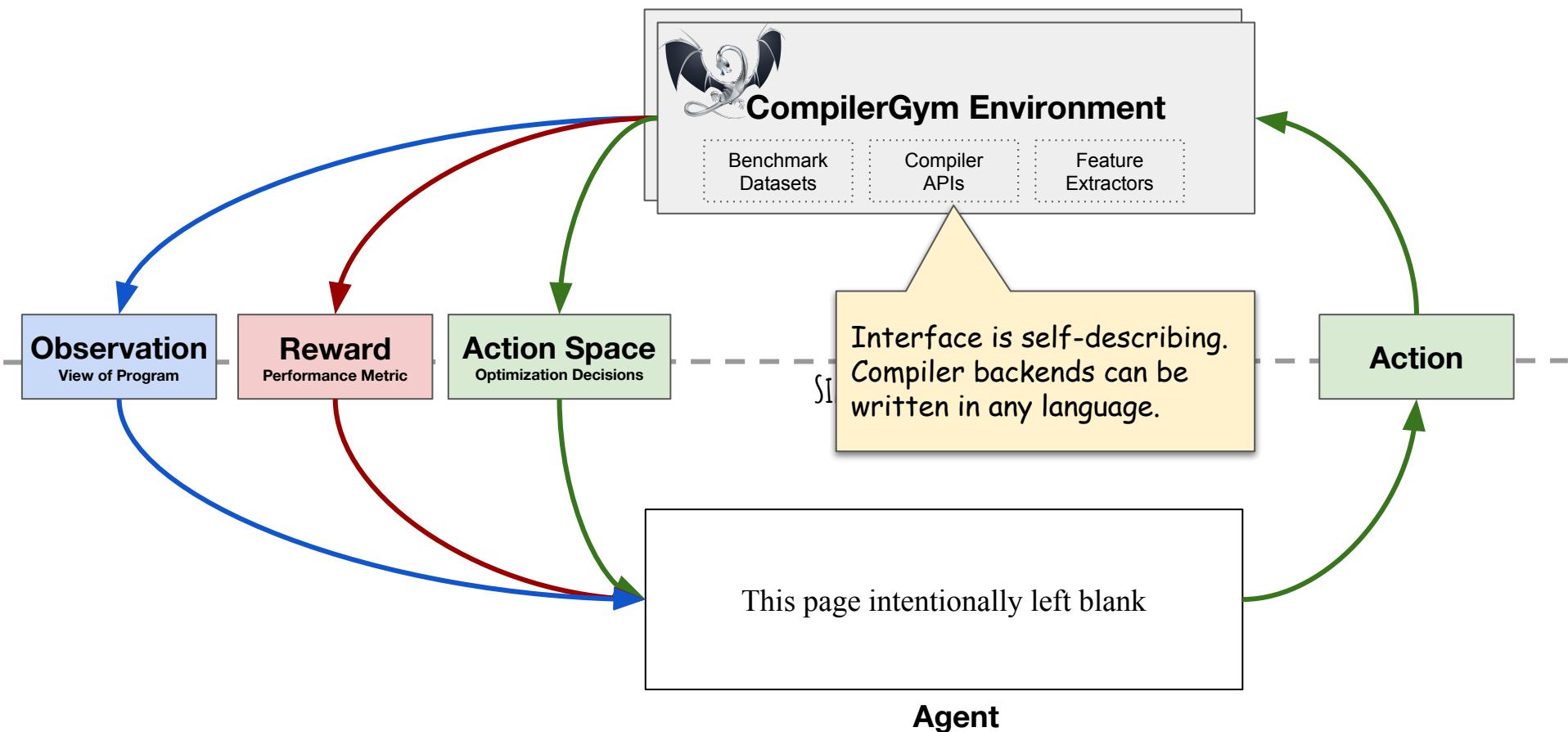
CompilerGym Goals

1. Lower the barrier to entry to AI for compilers
2. Advance the state-of-the-art in AI for compilers
3. Provide common benchmarks for compiler tasks









```
$ pip install compiler_gym
```

```
# no compilation or config
```

```
$ pip install compiler_gym          # no compilation or config  
$ python  
>>> import gym  
>>> import compiler_gym            # ships with compiler environments
```

```
$ pip install compiler_gym          # no compilation or config
$ python
>>> import gym
>>> import compiler_gym            # ships with compiler environments
>>> env = gym.make(
...     "llvm-v0",                  # creates a new environment
...     benchmark="cbench-v1/qsort",  # selects a compiler to use
...     observation_space="Autophase",# the program to compile
...     reward_space="IrInstructionCount0z", # the observation space
... )                                # the optimization target
```

```
$ pip install compiler_gym          # no compilation or config
$ python
>>> import gym
>>> import compiler_gym            # ships with compiler environments
>>> env = gym.make(
...     "llvm-v0",                  # creates a new environment
...     benchmark="cbench-v1/qsort",  # selects a compiler to use
...     observation_space="Autophase", # the program to compile
...     reward_space="IrInstructionCount0z", # the observation space
...     )                            # the optimization target
>>> observation = env.reset(benchmark)
>>> for _ in range(1000):
...     action = env.action_space.sample()    # your idea here!
...     observation, reward, done, info = env.step(action)
...     if done:
...         observation = env.reset(benchmark)
```

```
$ pip install compiler_gym          # no compilation or config
$ python
>>> import gym
>>> import compiler_gym            # ships with compiler environments
>>> env = gym.make(
...     "llvm-v0",                  # creates a new environment
...     benchmark="cbench-v1/qsort",  # selects a compiler to use
...     observation_space="Autophase", # the program to compile
...     reward_space="IrInstructionCount0z", # the observation space
...     )                            # the optimization target
>>> observation = env.reset(benchmark)
>>> for _ in range(1000):
...     action = env.action_space.sample()    # your idea here!
...     observation, reward, done, info = env.step(action)
...     if done:
...         observation = env.reset(benchmark)
>>> env.write_bitcode("benchmark.bc")      # save program to disk for later analysis
```

CompilerGym Features

- Three compiler problems included:
 - LLVM phase ordering
 - CUDA code generation
 - GCC flag tuning
- Batteries included
 - **Millions of programs** for training, 41x more than prior works
 - Huge optimization spaces
 - Several observation spaces from published work
 - Optimizing for **reward** or **codesize**
- Battle tested
- **27x faster** than prior works

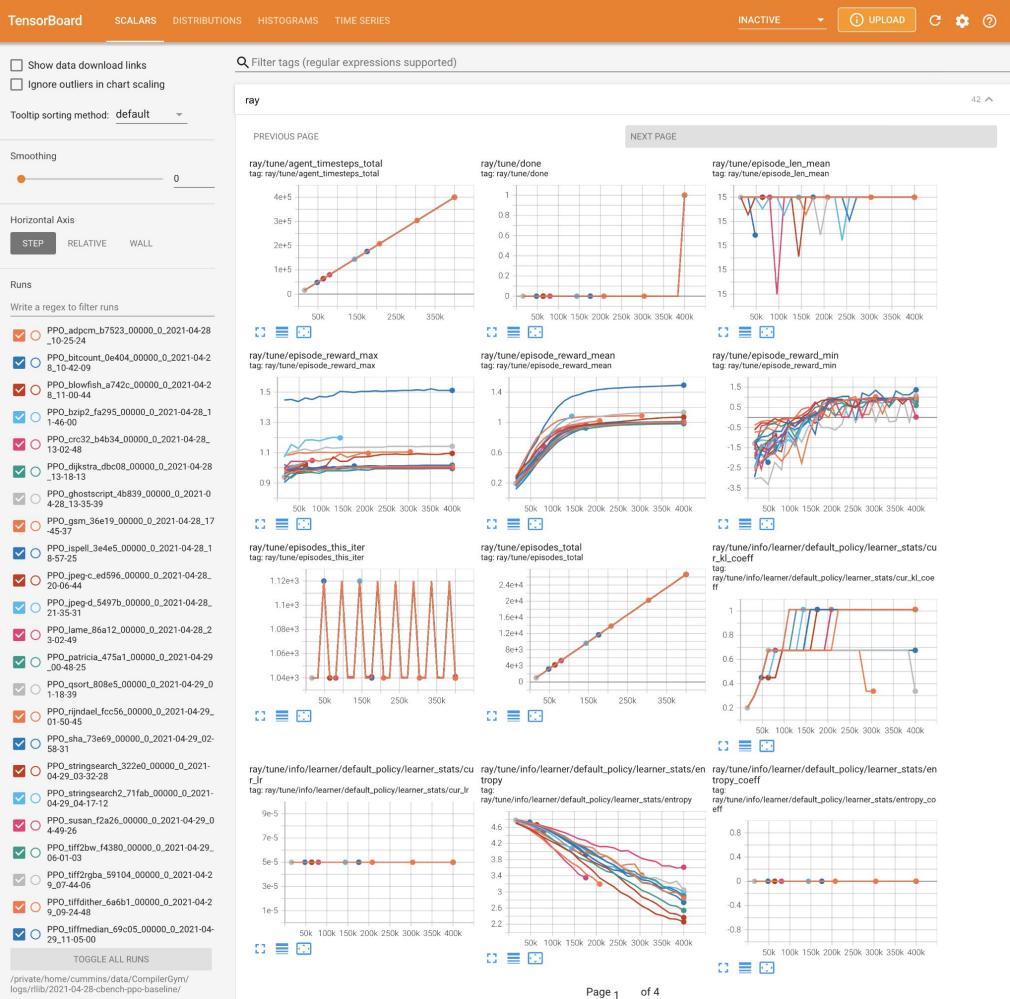
Why I think its cool

Repro state-of-the-art in < 15 LOC

```
import gym
import compiler_gym
from ray import tune
from ray.rllib.agents.ppo import PPOTrainer

tune.register_env(
    "CompilerGym",
    lambda _: gym.wrappers.TimeLimit(
        max_episode_steps=45,
        env=compiler_gym.make(
            "llvm-autophase-ic-v0",
            benchmark="cbench-v1/qsort",
        ),
    ),
)

tune.run(PPOTrainer, config={"env": "CompilerGym"})
```



Public Leaderboards



Author	Algorithm	Date	Walltime (lower is better)	Code size reduction (higher is better)
Facebook	Nevergrad (default optimizer)	2021-05	3,545.635s	1.074x
Facebook	Random agent (t=10800)	2021-03	10,512.356s	1.062x
Facebook	Greedy Search	2021-03	169.237s	1.055x
Facebook	Tabular Q	2021-04	2,534.305	1.036x
Facebook	Random agent (t=10)	2021-03	42.939s	1.031x
Your name here	Your technique here			

Encourage friendly competition across common benchmarks.
Provide a platform for promoting (and validating) people's work.
Low effort submission via pull requests.

Making Compiler Optimization Accessible to All

1. **Compiler research is fun** 😊

- Challenging, {high-dimensional, noisy} {action, state, reward} spaces
- Impact! Lots of room to innovate. Publications and potential for real world use

2. **CompilerGym reduces barrier to entry to "pip install."**

- Lets AI researchers focus on solving the problems
- Lets compiler developers focus on the compiler

3. **Fast pace of development, lots to look forward to!**

- New compilers + runtime optimization for LLVM
- Collaborations with academia to push state-of-art

overview

1. What I used to do
2. What I now do
3. What you should do

Deep Learning for Compilers, Fun, & Profit

- **A PhD is a lot of fun!**
- **My research is on using ML to improve compilers.**
- **Now working on democratizing the field.**

Chris Cummins
<https://chriscummins.cc>
cummins@fb.com