

# Synthesizing Benchmarks for Predictive Modeling

<http://chriscummins.cc/cgo17>



## Chris Cummins

University of Edinburgh



## Pavlos Petoumenos

University of Edinburgh



## Zheng Wang

Lancaster University



## Hugh Leather

University of Edinburgh

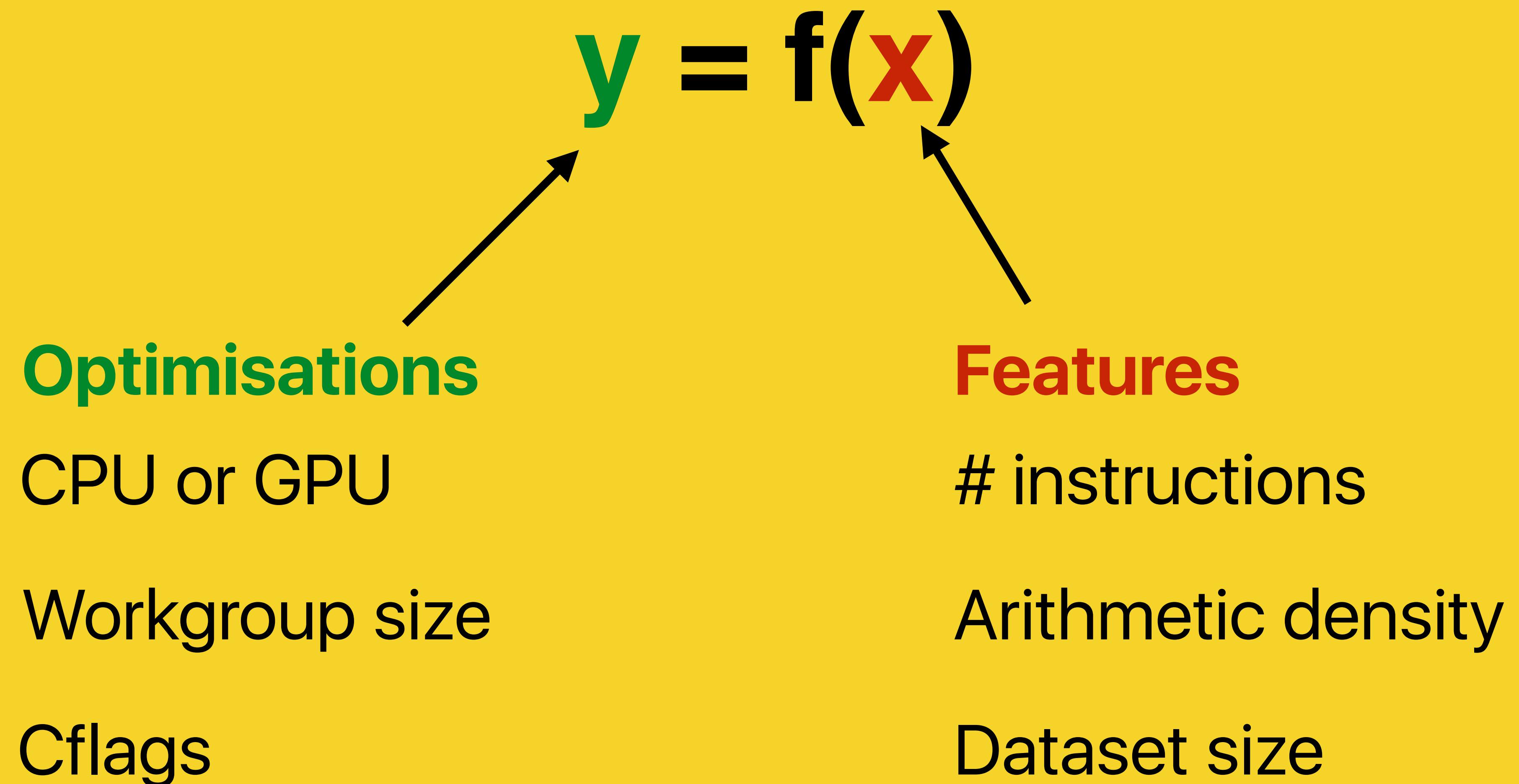


THE UNIVERSITY of EDINBURGH  
**informatics**

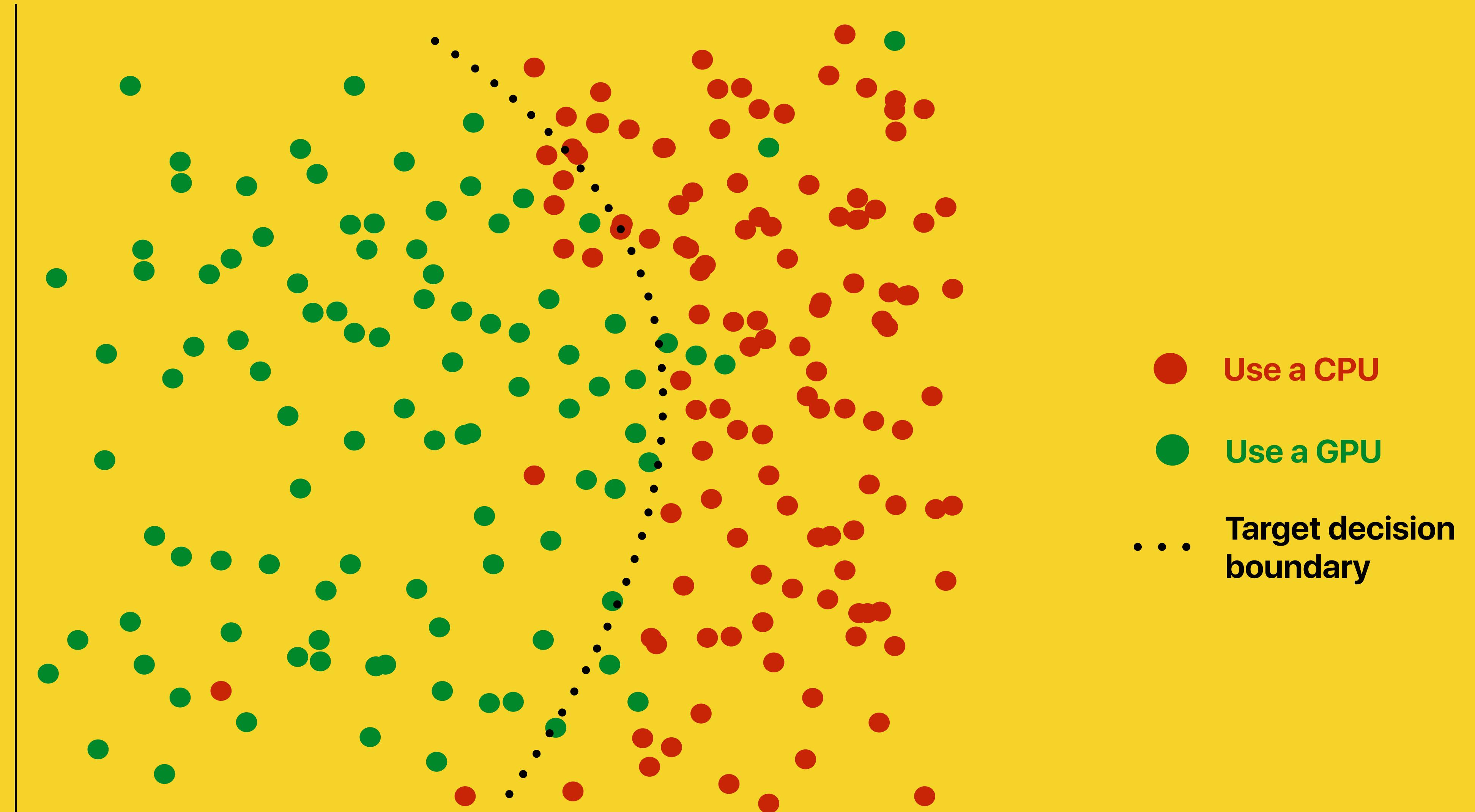
EPSRC Centre for Doctoral Training in  
**Pervasive Parallelism**

**EPSRC**  
Engineering and Physical Sciences  
Research Council

# *machine learning for compilers*

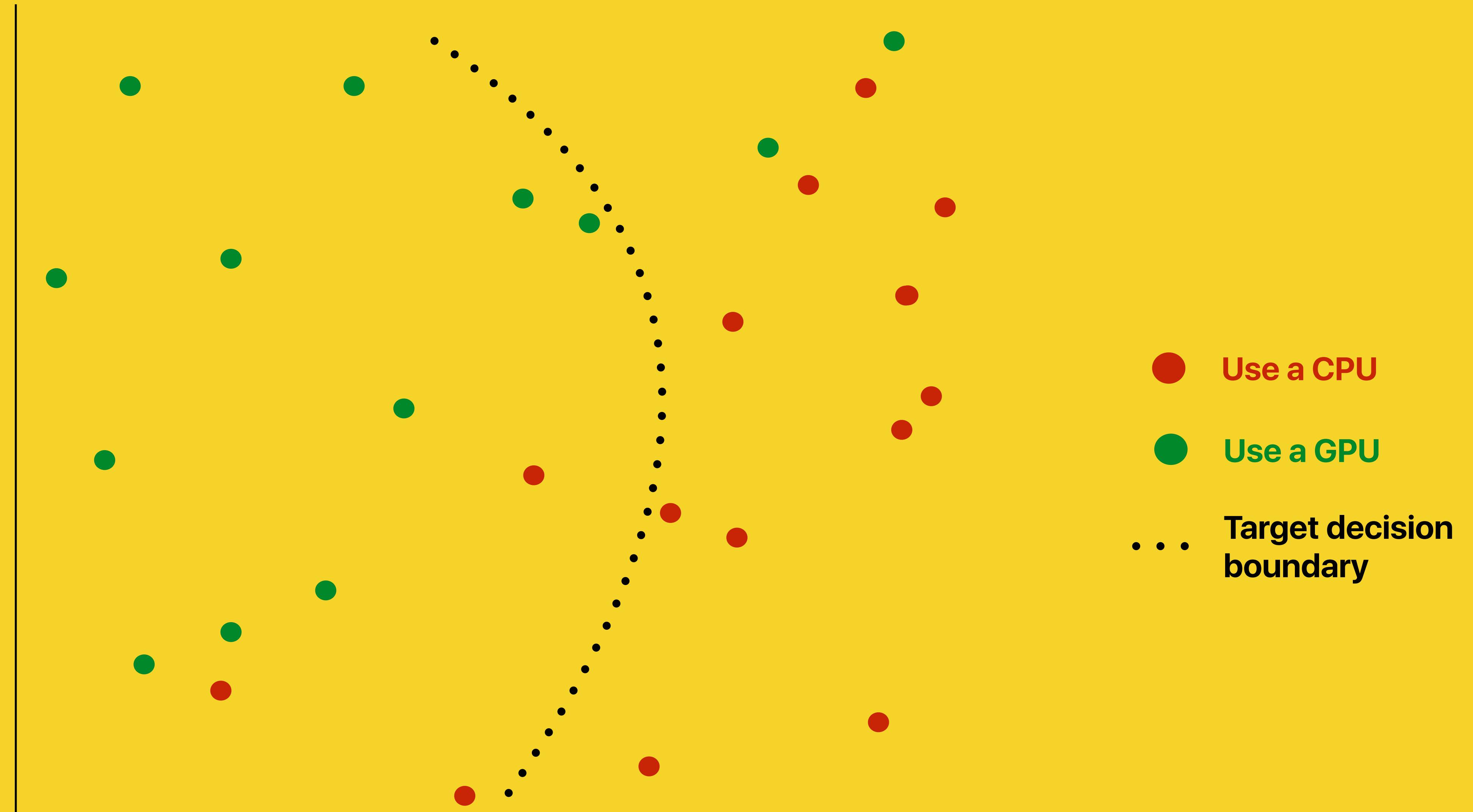


# *machine learning for compilers*



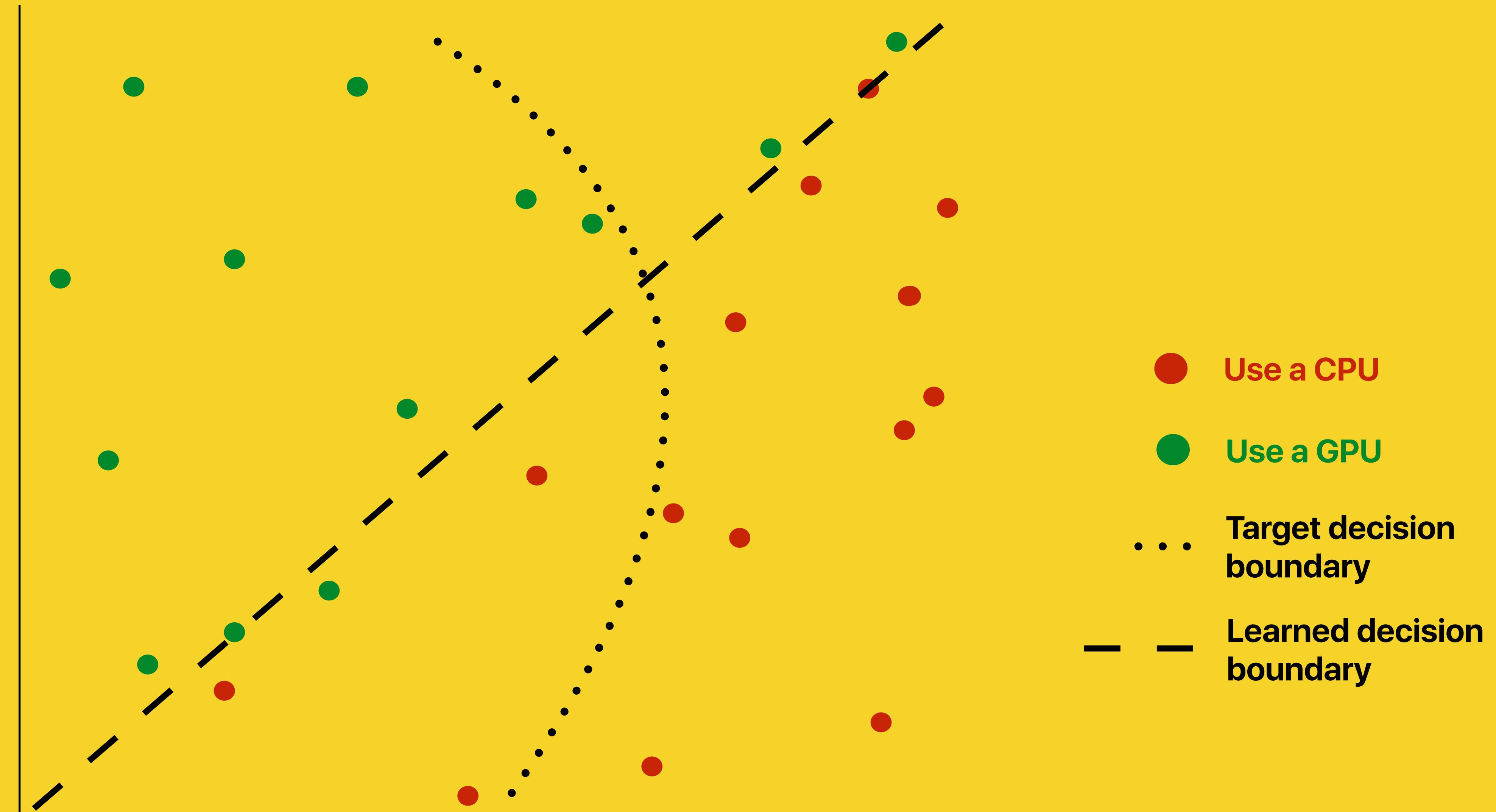
The idea.

# *machine learning for compilers*



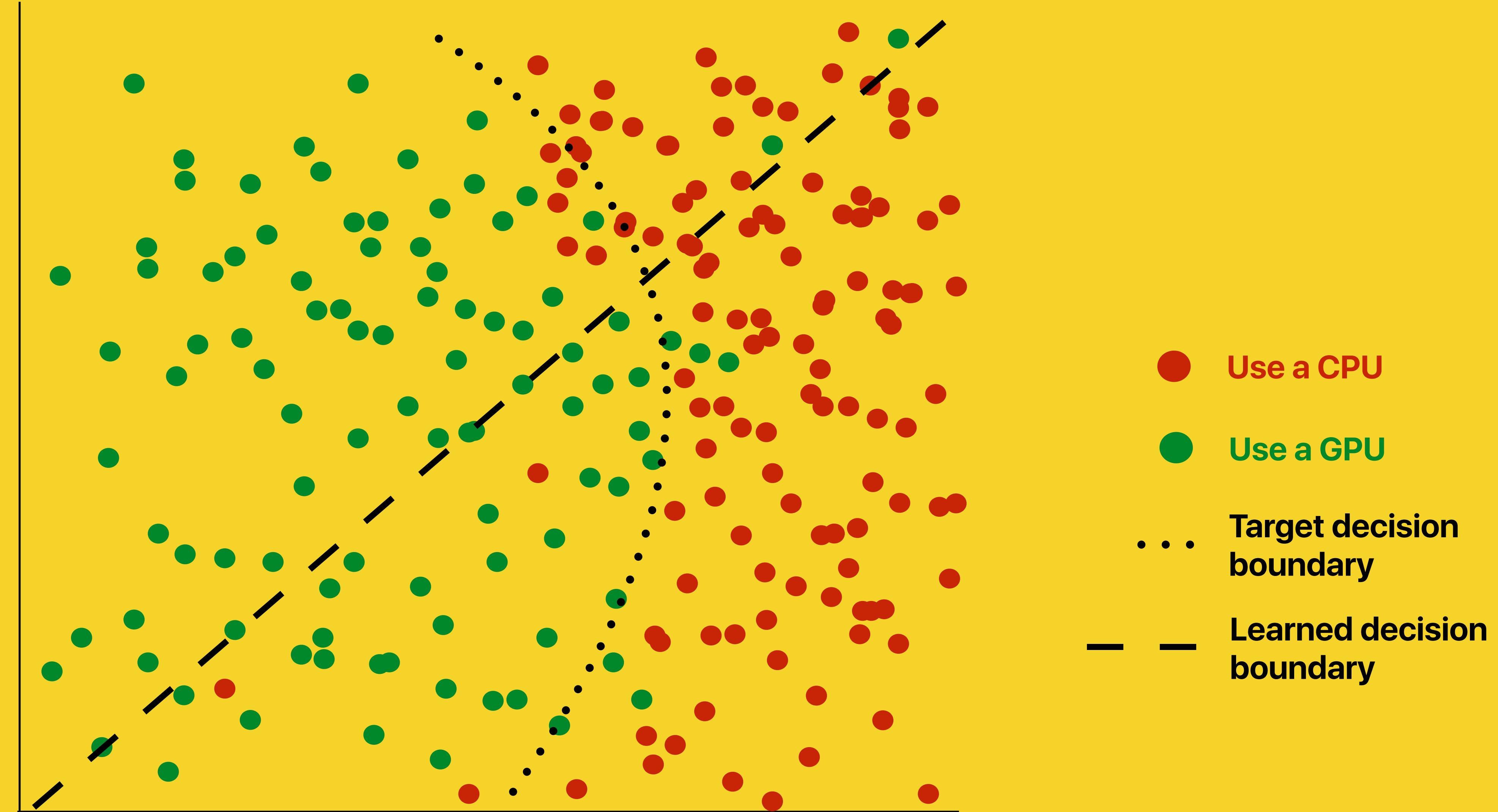
The reality.

# *machine learning for compilers*



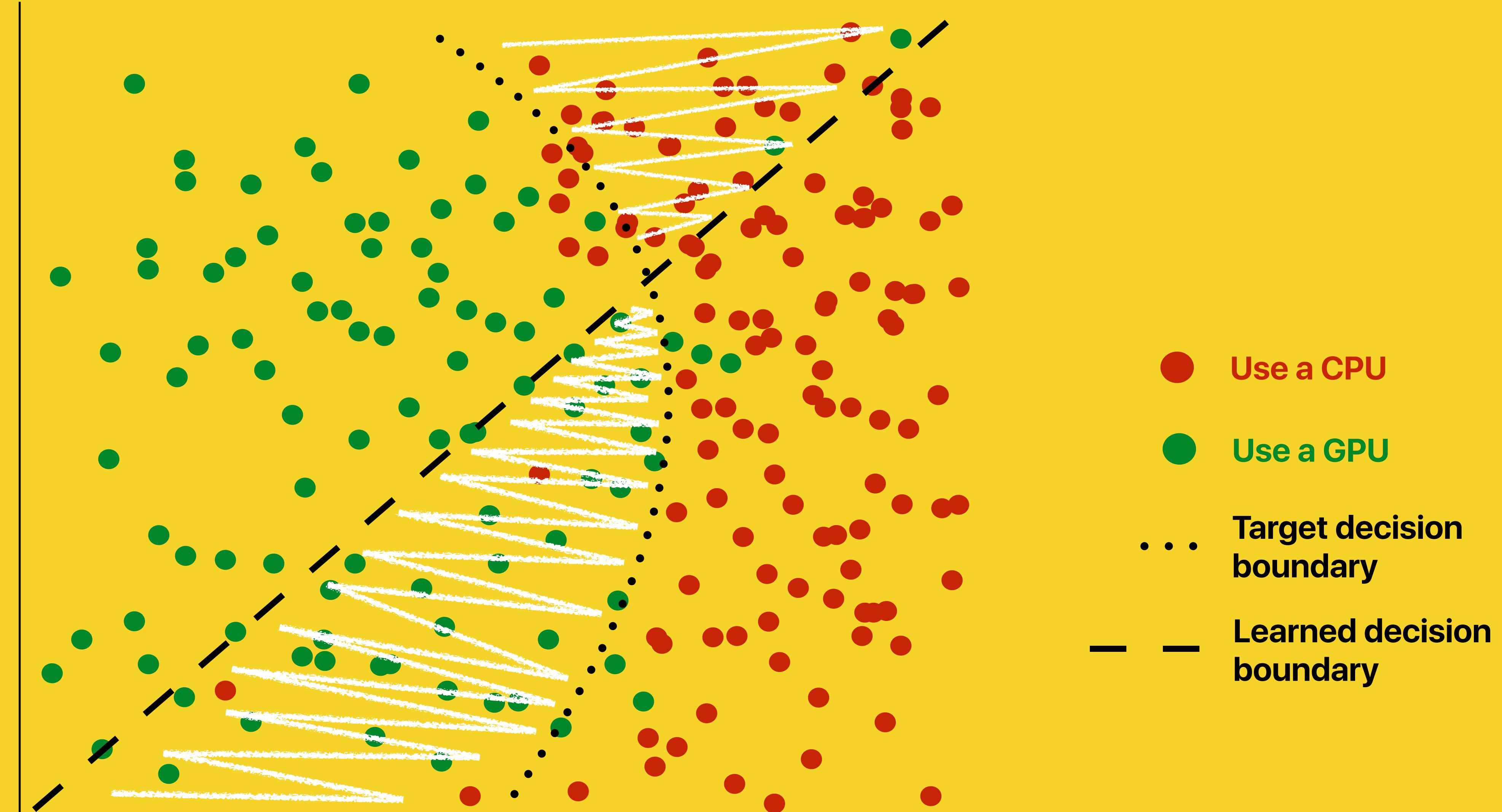
The reality.

# *machine learning for compilers*



**sparse data leads to inaccurate models!**

# *machine learning for compilers*



**sparse data leads to inaccurate models!**

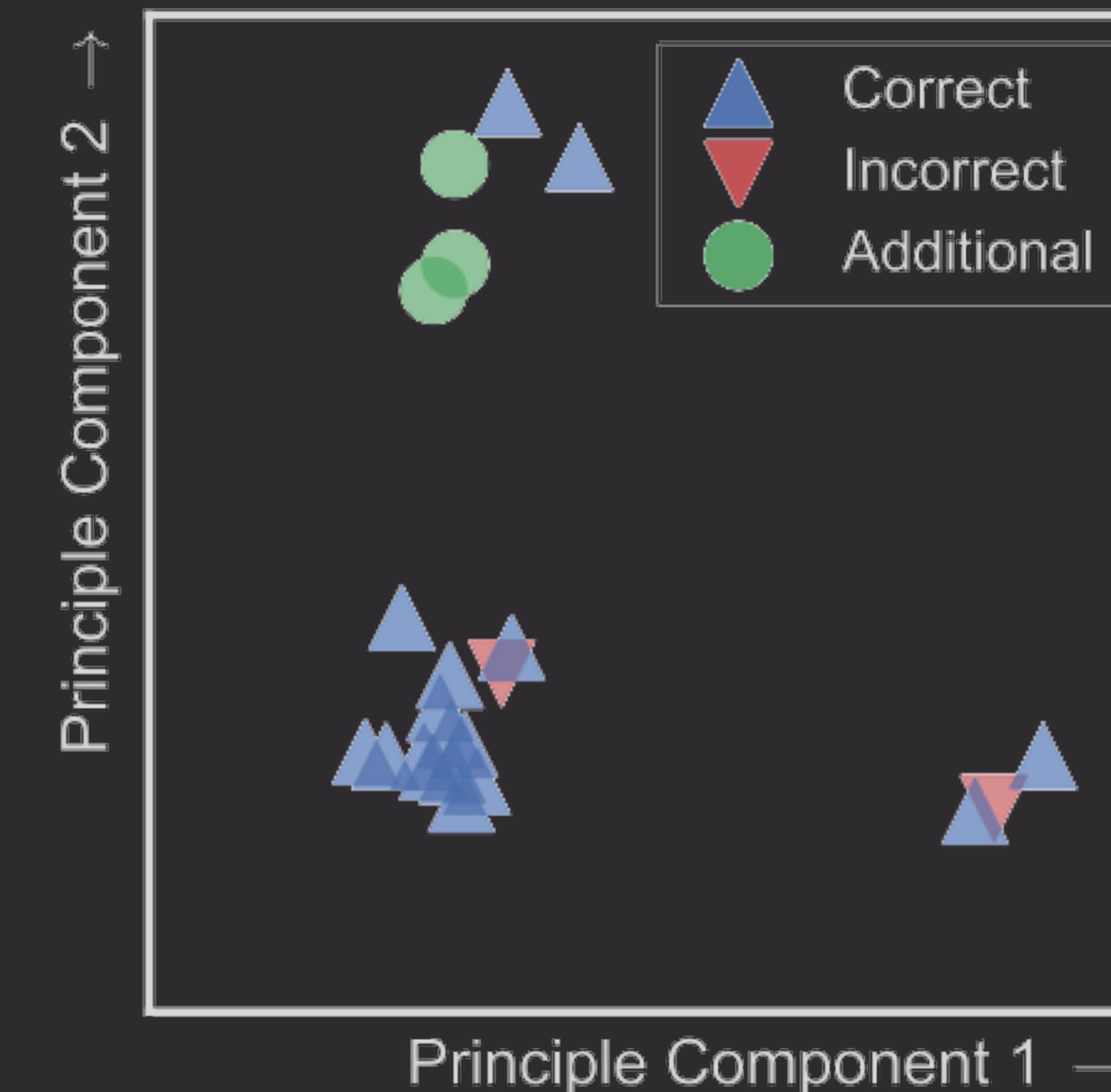
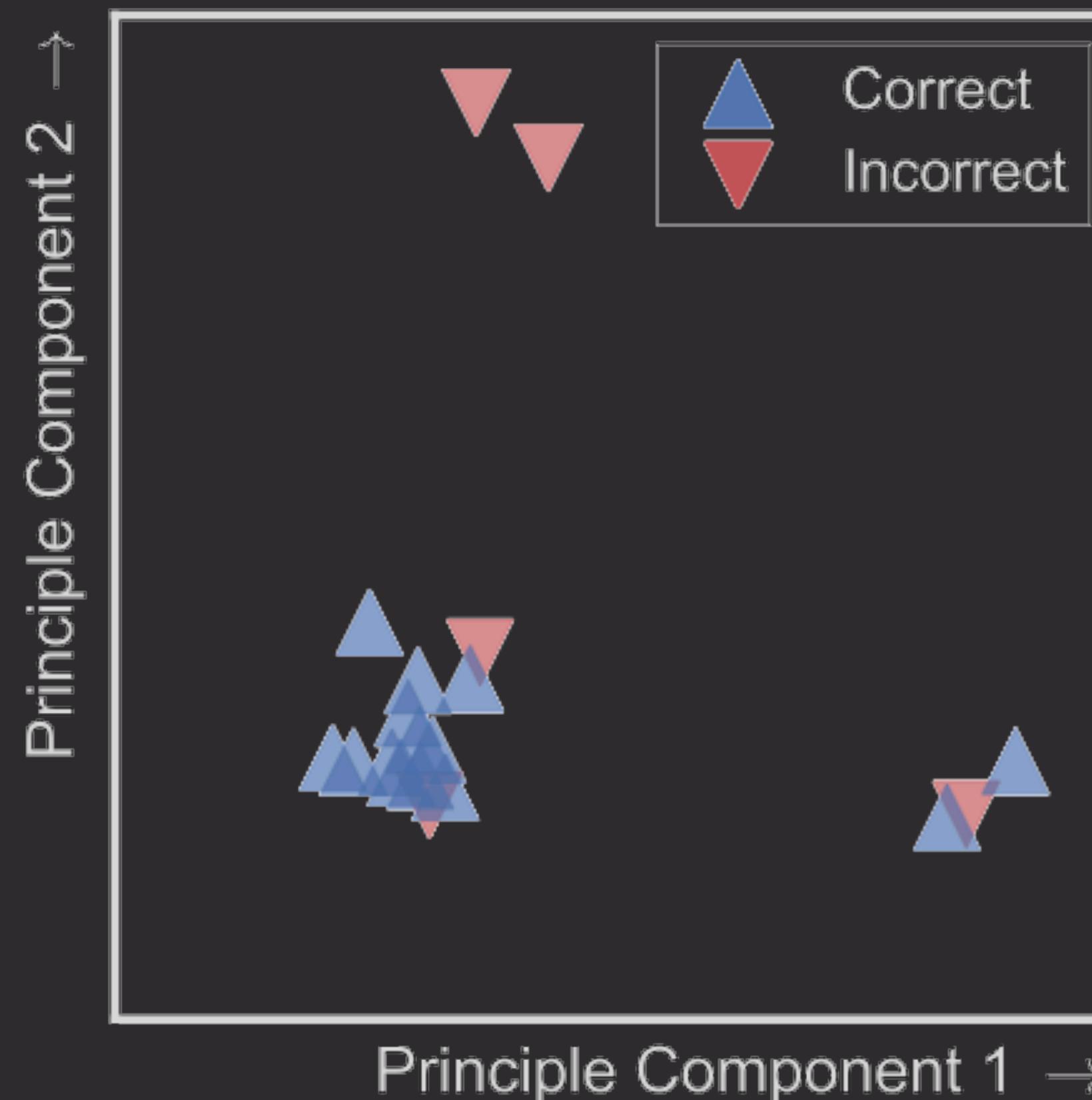
# *problem statement*

## 1. there aren't enough benchmarks

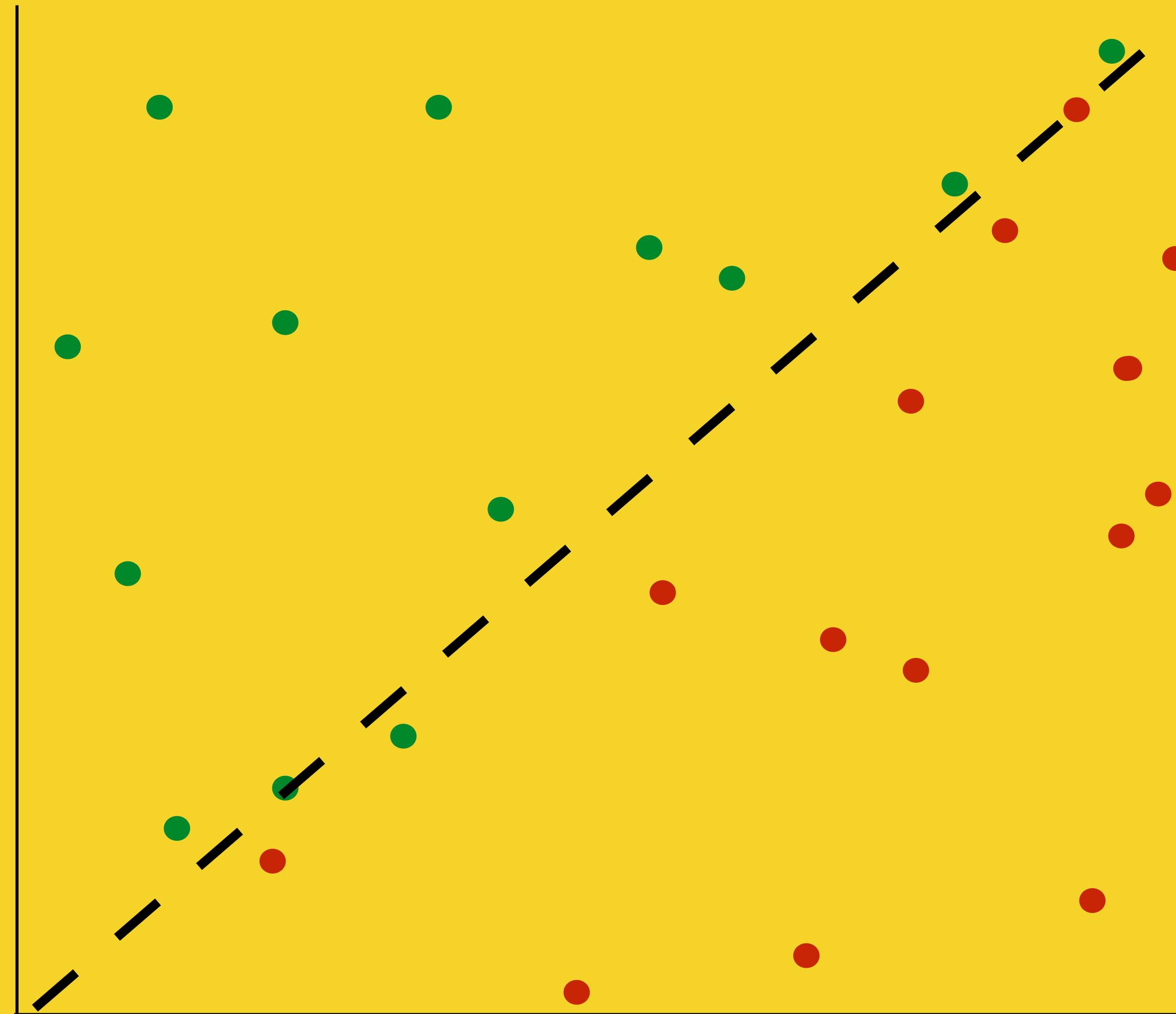
avg compiler paper	17	
Iris dataset	150	10x
MNIST dataset	60,000	$10^3$ x
ImageNet dataset	10,000,000	$10^6$ x

# *problem statement*

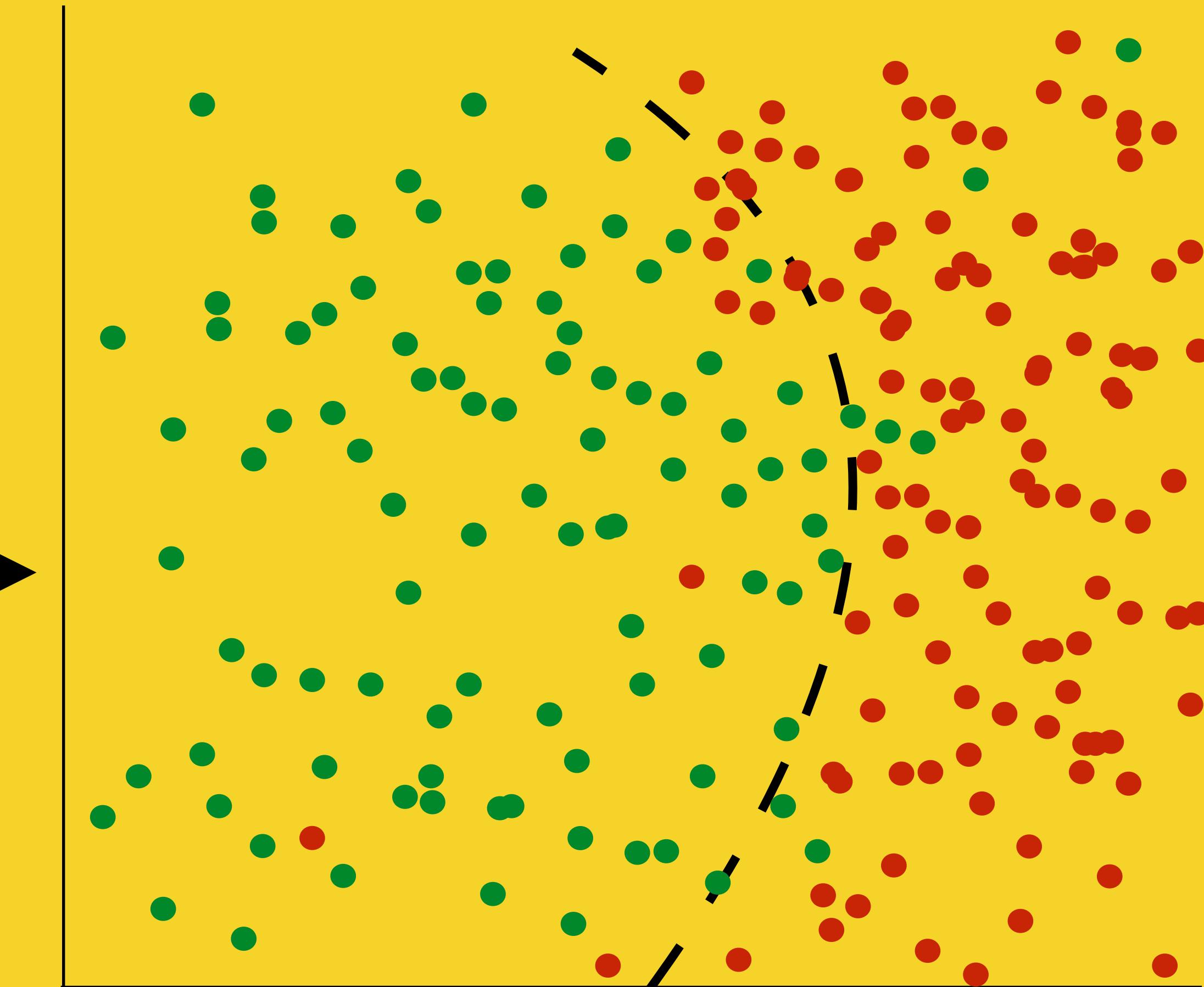
1. there aren't enough benchmarks
2. more benchmarks = better models



# *what we need*

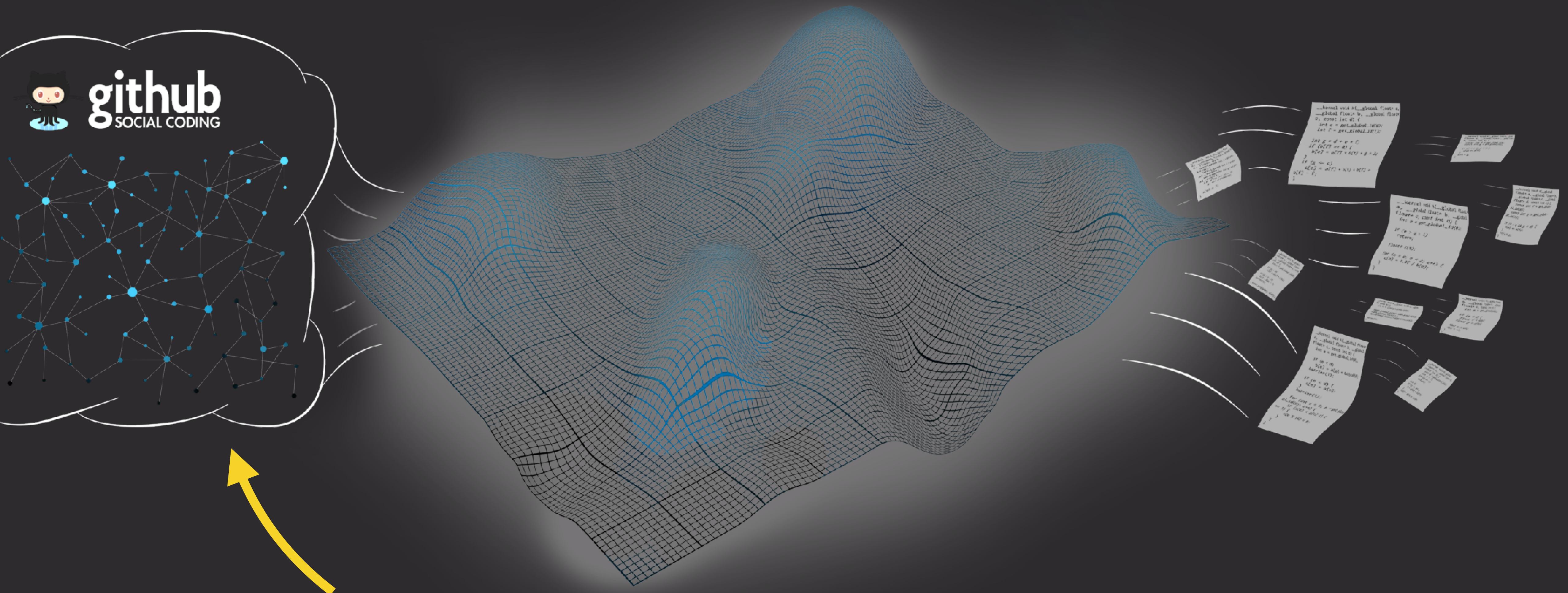


**from this**



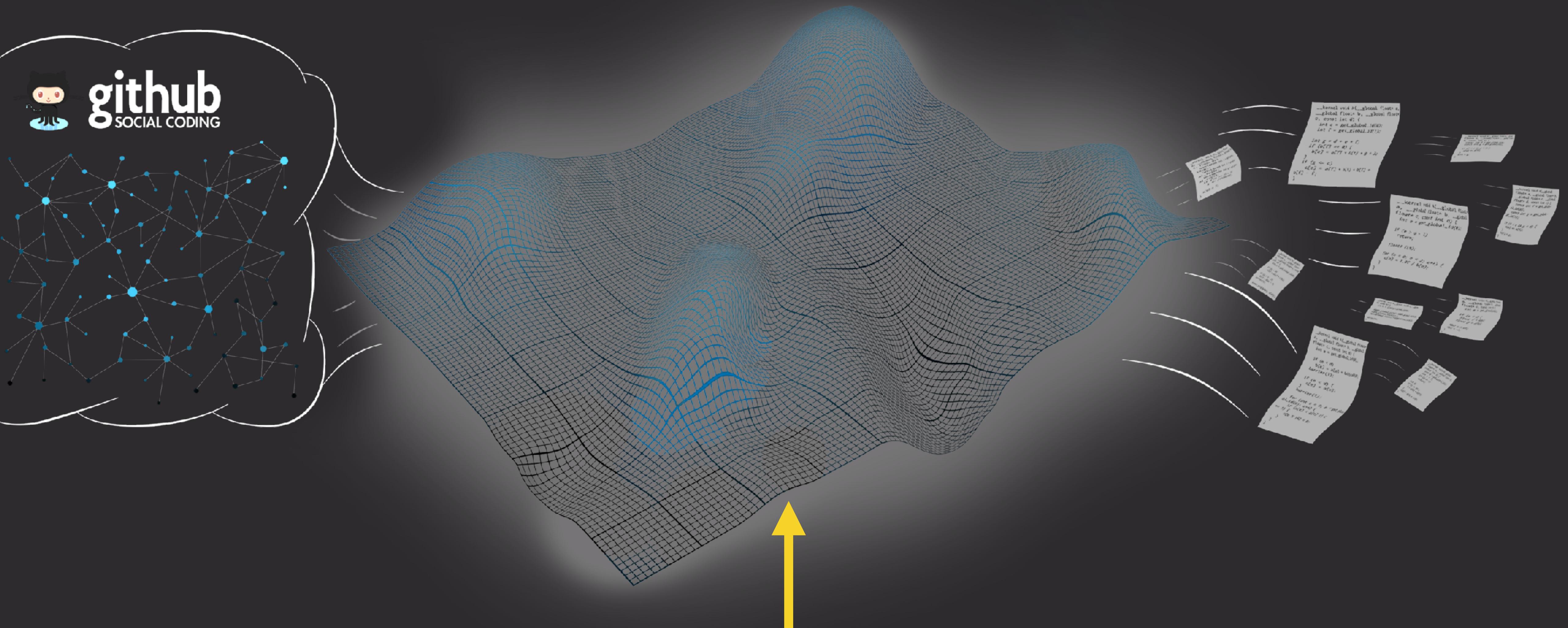
**to this**

# *our approach*



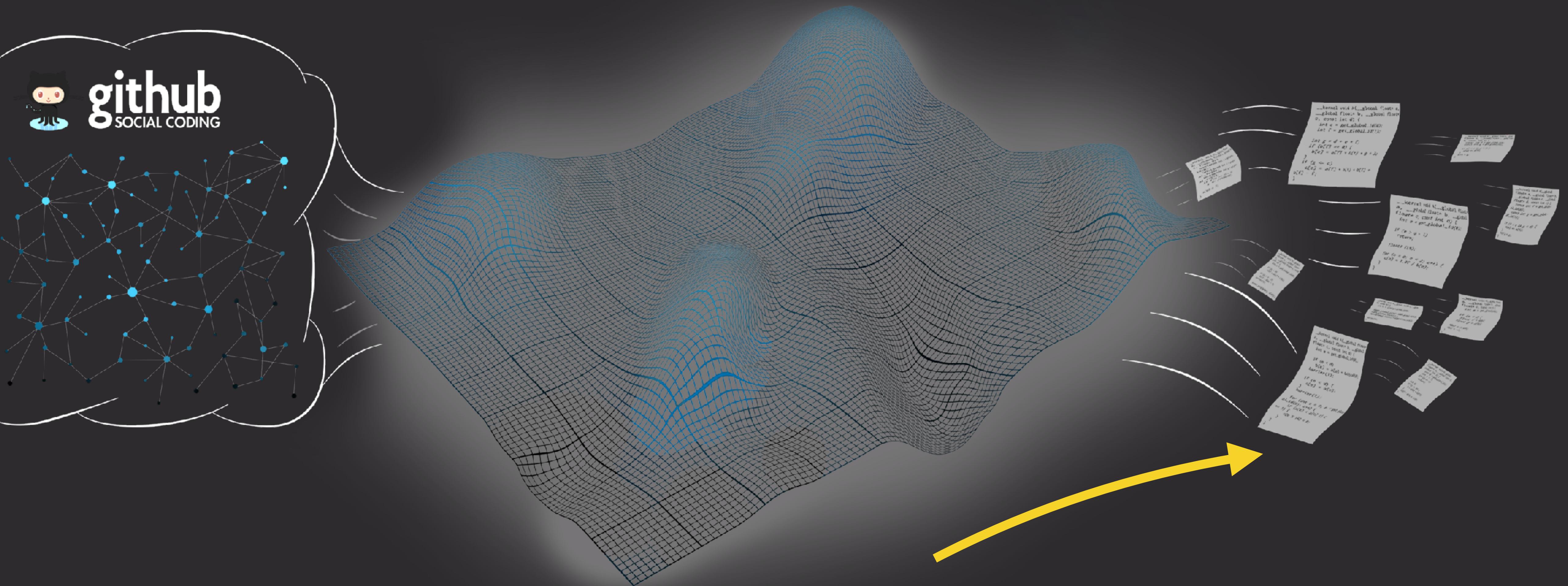
mine code from web

# *our approach*



# model source distr.

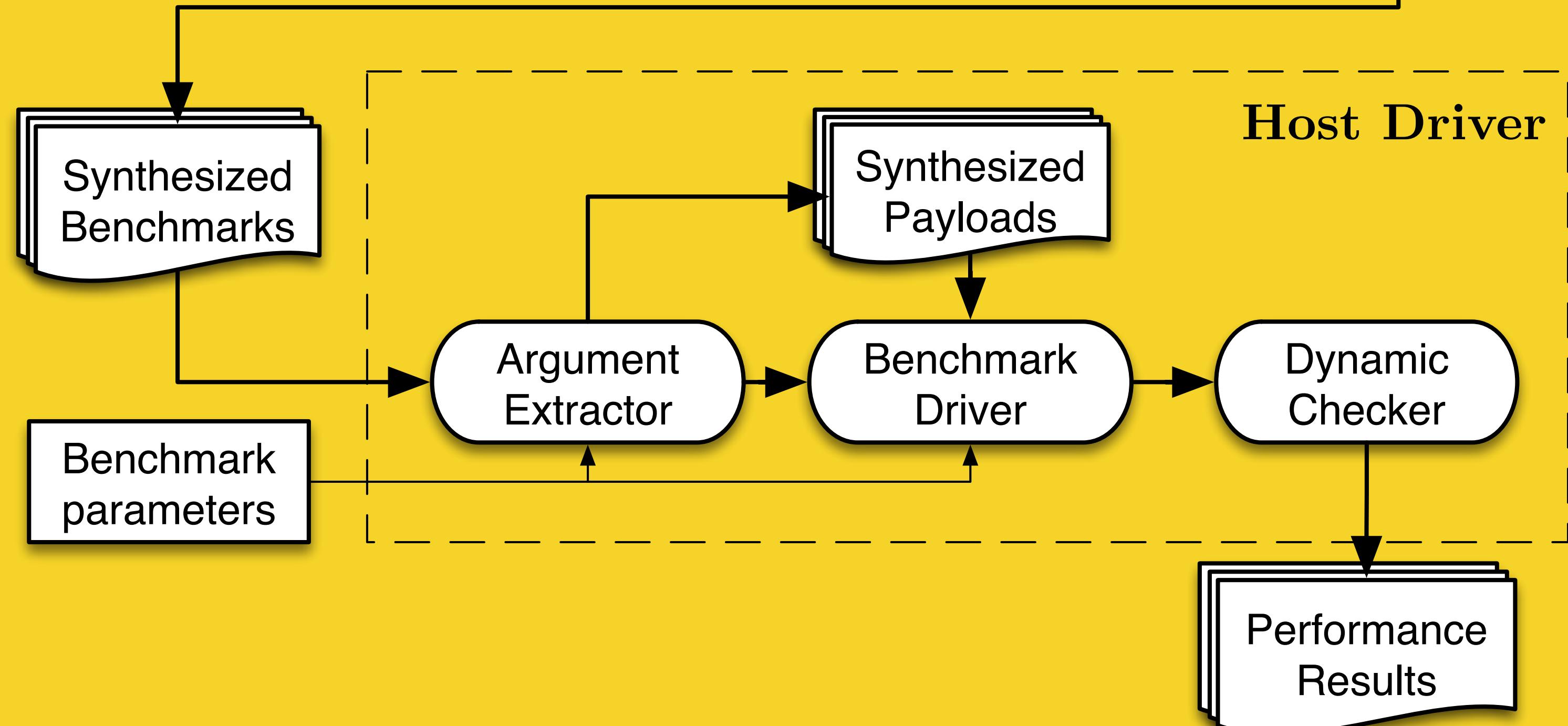
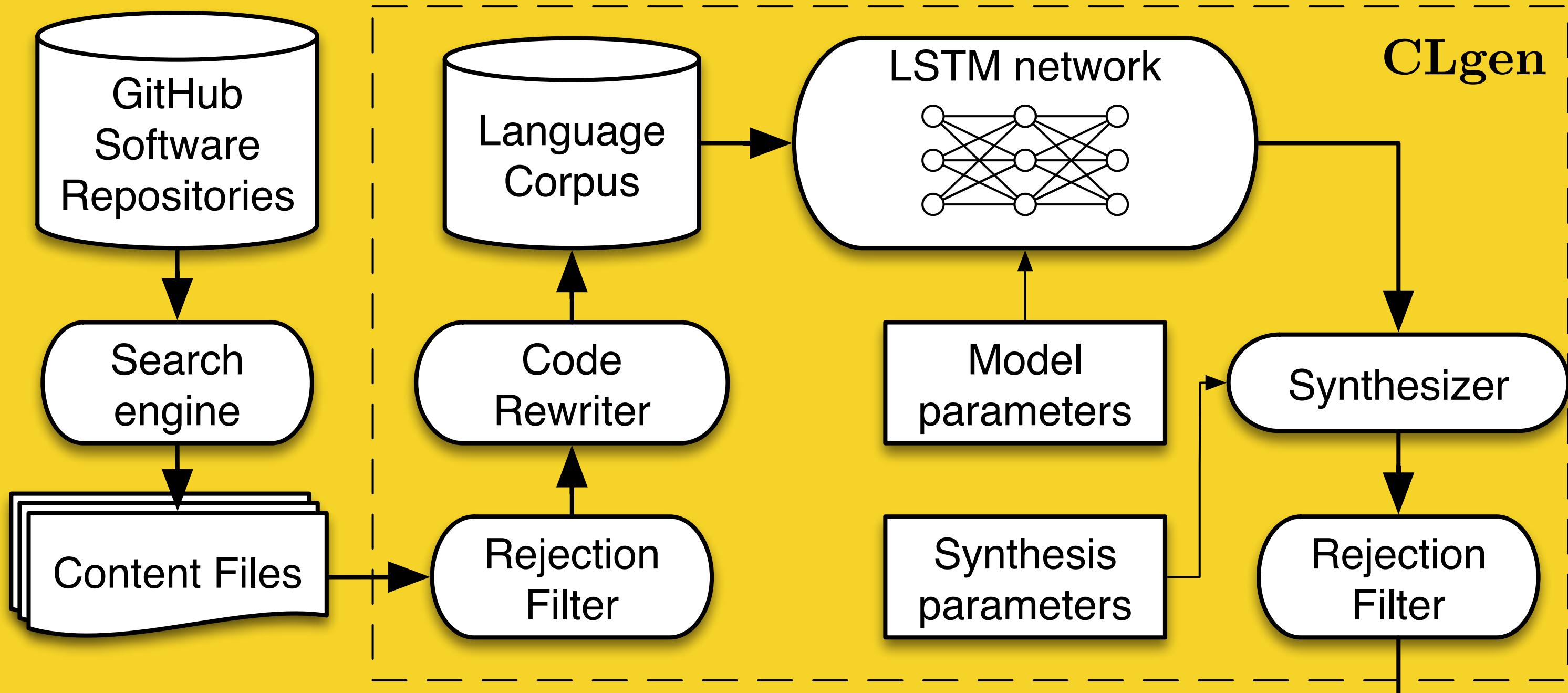
# *our approach*



# sample lang. model

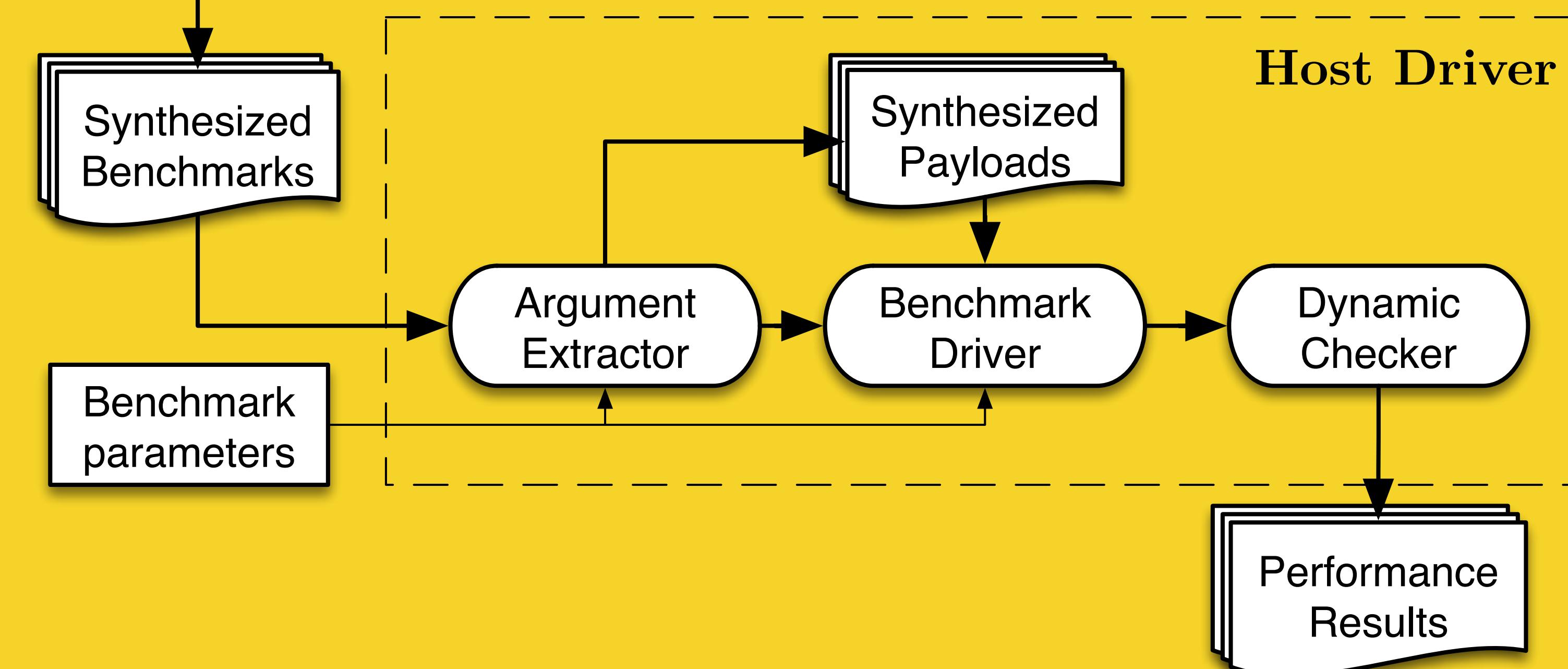
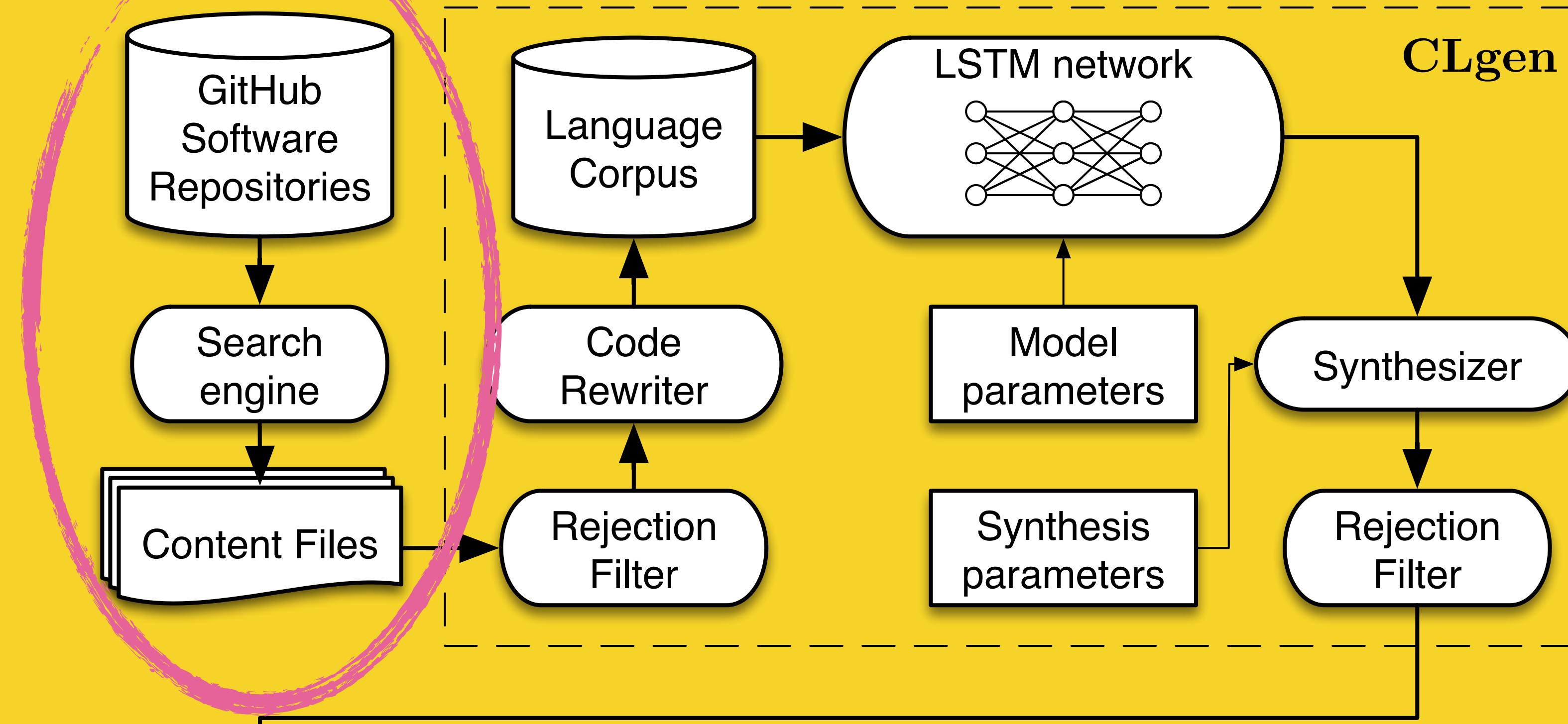
# *contributions*

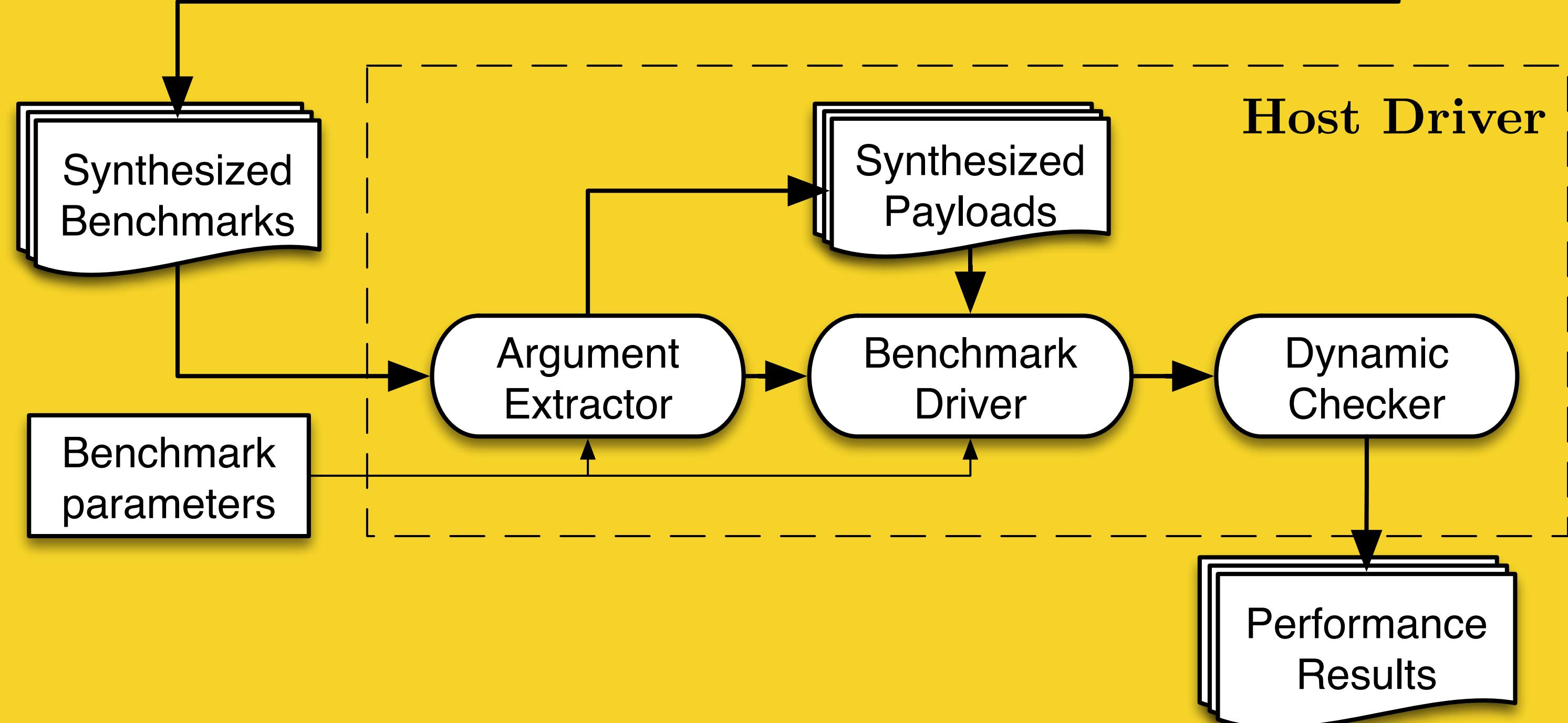
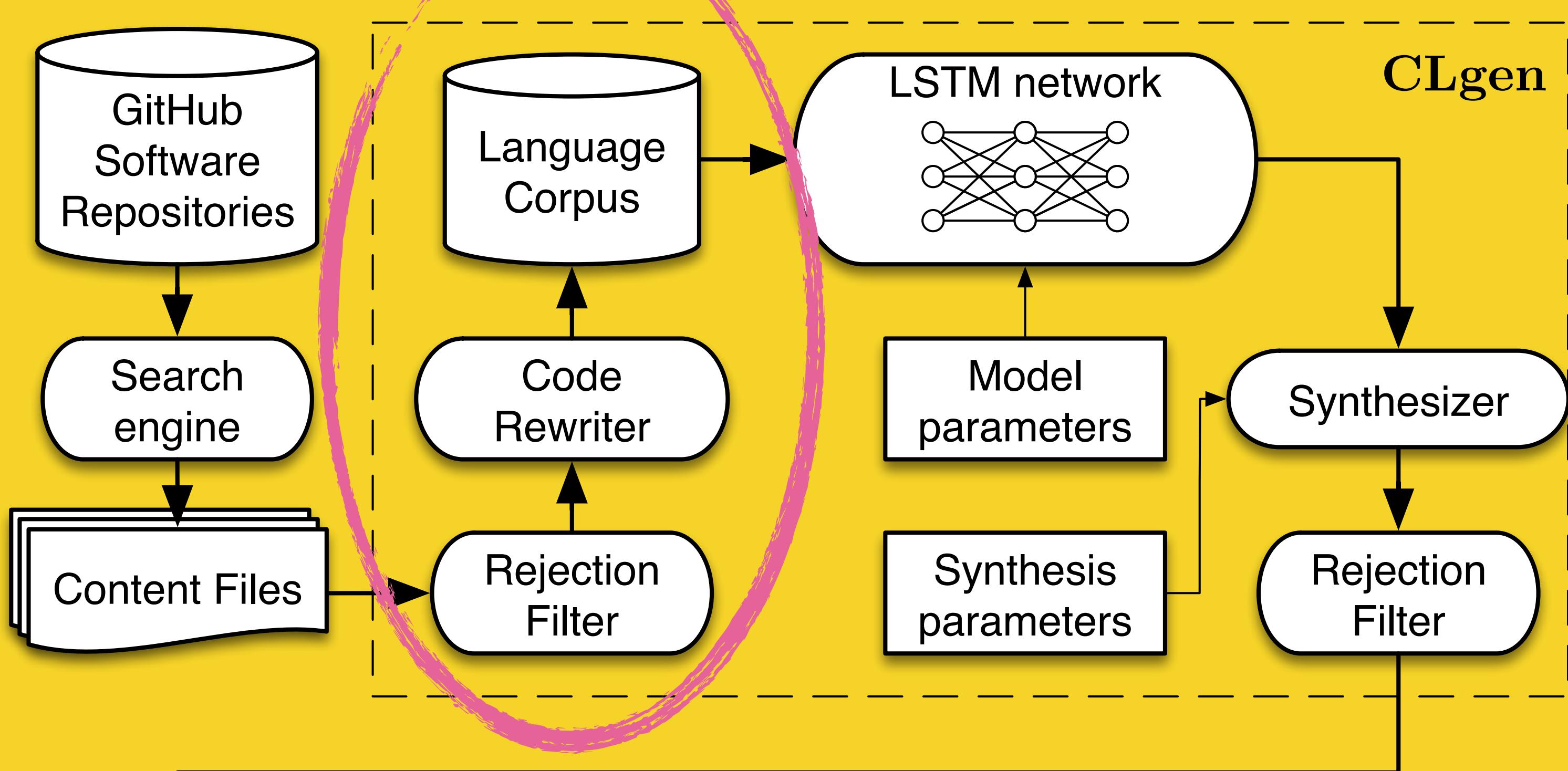
- 1. a deep learning approach to modeling PL semantics & usage**
- 2. first solution for general purpose benchmark synthesis**
- 3. automatic  $1.27\times$  speedup**
- 4. improved model designs for further  $4.30\times$  speedup**



8078  
files

2.8M  
lines





```
/* Copyright (C) 2014, Joe Blogs. */
#define CLAMPING
#define THRESHOLD_MAX 1.0f

float myclamp(float in) {
#ifdef CLAMPING
    return in > THRESHOLD_MAX ? THRESHOLD_MAX : in < 0.0f ? 0.0f : in;
#else
    return in;
#endif // CLAMPING
}

__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,
                                    int num_elems)
{
    // Do something really flipping cool
    int id = get_global_id(0);
    if (id < num_elems)
    {

        out[id] = myclamp(in[id]);
    }
}
```

```
/* Copyright (C) 2014, Joe Blogs. */
#define CLAMPING
#define THRESHOLD_MAX 1.0f

float myclamp(float in) {
#ifdef CLAMPING
    return in > THRESHOLD_MAX ? THRESHOLD_MAX : in < 0.0f ? 0.0f : in;
#else
    return in;
#endif // CLAMPING
}

__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,
                                    int num_elems)
{
    // Do something really flipping cool
    int id = get_global_id(0);
    if (id < num_elems)
    {
        out[id] = myclamp(in[id]);
    }
}
```

Is this real, valid OpenCL?  
Can we minimise non-functional variance?

*Strip comments*

```
/* Copyright (C) 2014, Joe Blogs. */  
#define CLAMPING  
#define THRESHOLD_MAX 1.0f  
  
float myclamp(float in) {  
#ifdef CLAMPING  
    return in > THRESHOLD_MAX ? THRESHOLD_MAX : in < 0.0f ? 0.0f : in;  
#else  
    return in;  
#endif // CLAMPING  
}  
  
kernel void findAllNodesMergedAabb(__global float* in, __global float* out,  
                                    int num_elems)  
{  
    // Do something really flipping cool  
    int id = get_global_id(0);  
    if (id < num_elems)  
    {  
        out[id] = myclamp(in[id]);  
    }  
}
```

~~/\* Copyright (C) 2011, Joe Blogs. \*/~~

#define CLAMPING  
#define THRESHOLD\_MAX 1.0f

float myclamp(float in) {  
#ifdef CLAMPING  
 return in > THRESHOLD\_MAX ? THRESHOLD\_MAX : in < 0.0f ? 0.0f : in;  
#else  
 return in;  
#endif // CLAMPING  
}

\_\_kernel void findAllNodesMergedAabb(\_\_global float\* in, \_\_global float\* out,  
int num\_elems)  
{  
// Do something really flipping cool  
int id = get\_global\_id(0);  
if (id < num\_elems)  
{  
  
 out[id] = myclamp(in[id]);  
}  
}

*Strip comments*

~~Strip comments~~

```
#define CLAMPING
#define THRESHOLD_MAX 1.0f

float myclamp(float in) {
#ifdef CLAMPING
    return in > THRESHOLD_MAX ? THRESHOLD_MAX : in < 0.0f ? 0.0f : in;
#else
    return in;
#endif}

__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,
                                    int num_elems)
{
    int id = get_global_id(0);
    if (id < num_elems)
    {

        out[id] = myclamp(in[id]);
    }
}
```

~~Strip comments~~  
Preprocess

```
#define CLAMPING
#define THRESHOLD_MAX 1.0f

float myclamp(float in) {
#ifdef CLAMPING
    return in > THRESHOLD_MAX ? THRESHOLD_MAX : in < 0.0f ? 0.0f : in;
#else
    return in;
#endif}

__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,
                                    int num_elems)
{
    int id = get_global_id(0);
    if (id < num_elems)
    {

        out[id] = myclamp(in[id]);
    }
}
```

~~Strip comments~~  
Preprocess

```
#define CLAMPING
#define THRESHOLD_MAX 1.0f

float myclamp(float in) {
#ifndef CLAMPING
    return in > THRESHOLD_MAX ? THRESHOLD_MAX : in < 0.0f ? 0.0f : in;
#else
    return in;
#endif}

__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,
                                    int num_elems)
{
    int id = get_global_id(0);
    if (id < num_elems)
    {

        out[id] = myclamp(in[id]);
    }
}
```

~~Strip comments~~  
~~Preprocess~~

```
float myclamp(float in) {
    return in > 1.0f ? 1.0f : in < 0.0f ? 0.0f : in;
}

__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,
                                    int num_elems)
{
    int id = get_global_id(0);
    if (id < num_elems)
    {

        out[id] = myclamp(in[id]);
    }
}
```

~~Strip comments~~  
~~Preprocess~~

```
float myclamp(float in) {
    return in > 1.0f ? 1.0f : in < 0.0f ? 0.0f : in;
}

__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,
                                    int num_elems)
{
    int id = get_global_id(0);
    if (id < num_elems)
    {

        out[id] = myclamp(in[id]);
    }
}
```

Does it compile?  
Does it contain instructions?

~~Strip comments~~  
~~Preprocess~~

```
float myclamp(float in) {  
    return in > 1.0f ? 1.0f : in < 0.0f ? 0.0f : in;  
}
```

```
__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,  
                                     int num_elems)
```

```
{  
    int id = get_global_id(0);  
    if (id < num_elems)  
    {  
  
        out[id] = myclamp(in[id]);  
    }  
}
```

~~Does it compile?~~

~~Does it contain instructions?~~

(33% discard rate)

```
float myclamp(float in) {  
    return in > 1.0f ? 1.0f : in < 0.0f ? 0.0f : in;
```

```
}
```

```
__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,  
                                     int num_elems)
```

```
{
```

```
    int id = get_global_id(0);
```

```
    if (id < num_elems)
```

```
{
```

```
        out[id] = myclamp(in[id]);
```

```
}
```

```
}
```

~~Strip comments  
Preprocess  
Rewrite function names~~

~~Does it compile?~~

~~Does it contain instructions?~~

(33% discard rate)

```
float myclamp(float in) {
    return in > 1.0f ? 1.0f : in < 0.0f ? 0.0f : in;
}

__kernel void findAllNodesMergedAabb(__global float* in, __global float* out,
                                      int num_elems)
{
    int id = get_global_id(0);
    if (id < num_elems)
    {

        out[id] = myclamp(in[id]);
    }
}
```

~~Strip comments  
Preprocess  
Rewrite function names~~

~~Does it compile?~~

~~Does it contain instructions?~~

(33% discard rate)

```
float A(float in) {  
    return in > 1.0f ? 1.0f : in < 0.0f ? 0.0f : in;
```

```
}
```

```
__kernel void B(__global float* in, __global float* out,  
                int num_elems)
```

```
{
```

```
    int id = get_global_id(0);
```

```
    if (id < num_elems)
```

```
{
```

```
        out[id] = A(in[id]);
```

```
}
```

```
}
```

~~Strip comments  
Preprocess  
Rewrite function names~~

~~Does it compile?~~

~~Does it contain instructions?~~

(33% discard rate)

```
float A(float in) {
    return in > 1.0f ? 1.0f : in < 0.0f ? 0.0f : in;
}

__kernel void B(__global float* in, __global float* out,
                int num_elems)
{
    int id = get_global_id(0);
    if (id < num_elems)
    {

        out[id] = A(in[id]);
    }
}
```

~~Strip comments~~  
~~Preprocess~~  
~~Rewrite function names~~  
~~Rewrite variable names~~

~~Does it compile?~~

~~Does it contain instructions?~~

(33% discard rate)

```
float A(float in) {  
    return in > 1.0f ? 1.0f : in < 0.0f ? 0.0f : in;  
}  
  
__kernel void B(__global float* in, __global float* out,  
                int num_elems)  
{  
    int id = get_global_id(0);  
    if (id < num_elems)  
    {  
        out[id] = A(in[id]);  
    }  
}
```

Strip comments  
Preprocess  
Rewrite function names  
Rewrite variable names

Does it compile?  
Does it contain instructions?  
(33% discard rate)

```
float A(float a) {  
    return a > 1.0f ? 1.0f : a < 0.0f ? 0.0f : a;  
}
```

~~Strip comments~~  
~~Preprocess~~  
~~Rewrite function names~~  
~~Rewrite variable names~~

```
kernel void B(global float* b, global float* c,  
                int d)  
{  
    int e = get_global_id(0);  
    if (e < d)  
    {  
        c[e] = A(b[e]);  
    }  
}
```

~~Does it compile?~~

~~Does it contain instructions?~~

(33% discard rate)

```
float A(float a) {  
    return a > 1.0f ? 1.0f : a < 0.0f ? 0.0f : a;  
}  
  
_kernel void B(__global float* b, __global float* c,  
int d)  
  
{ int e = get_global_id(0);  
↳ if (e < d)  
{  
↳ X c[e] = A(b[e]);  
}  
}
```

~~Strip comments~~  
~~Preprocess~~  
~~Rewrite function names~~  
~~Rewrite variable names~~  
~~Enforce code style~~

~~Does it compile?~~

~~Does it contain instructions?~~

(33% discard rate)

```
float A(float a) {  
    return a > 1.0f ? 1.0f : a < 0.0f ? 0.0f : a;  
}
```

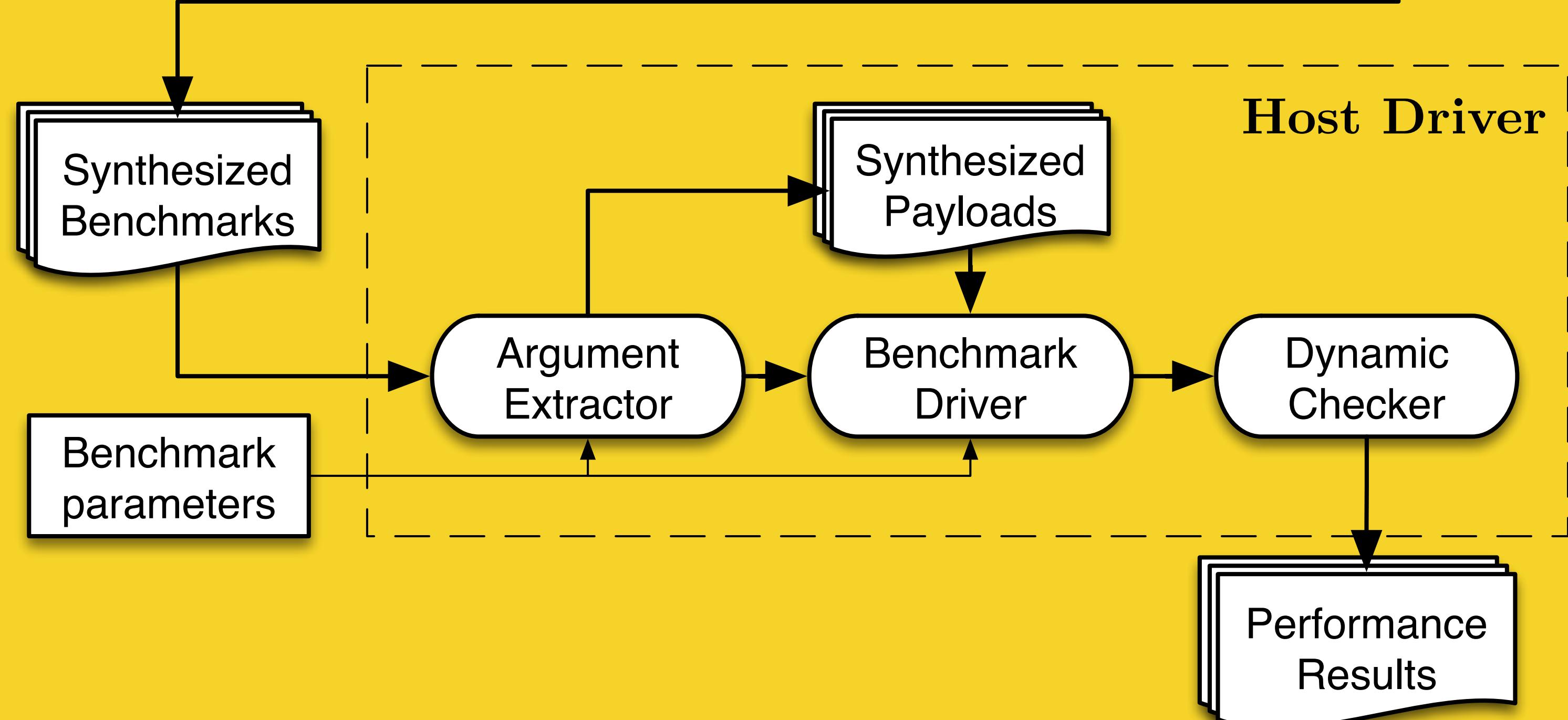
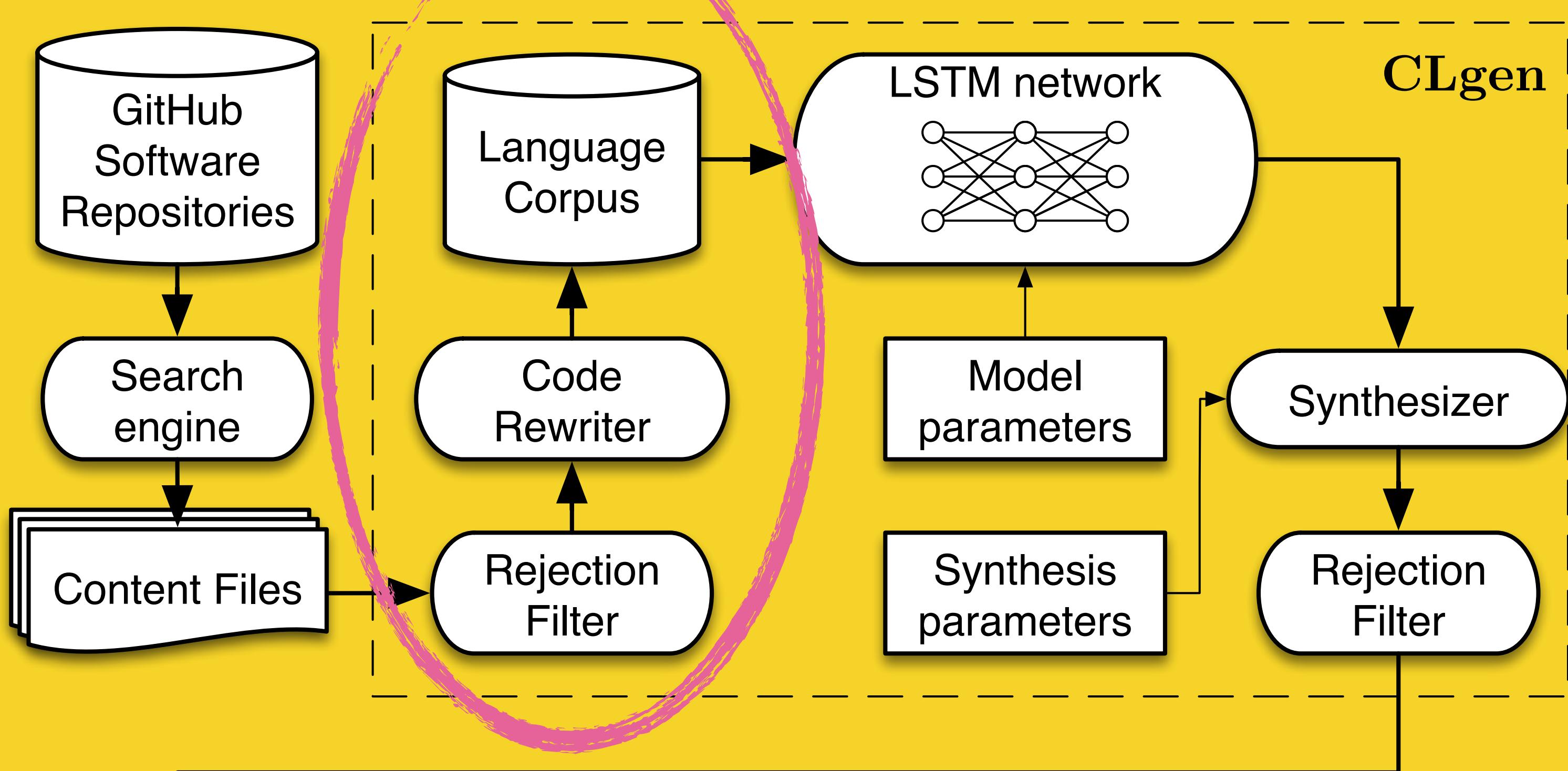
```
__kernel void B(__global float* b, __global float* c, int d) {  
    int e = get_global_id(0);  
    if (e < d) {  
        c[e] = A(b[e]);  
    }  
}
```

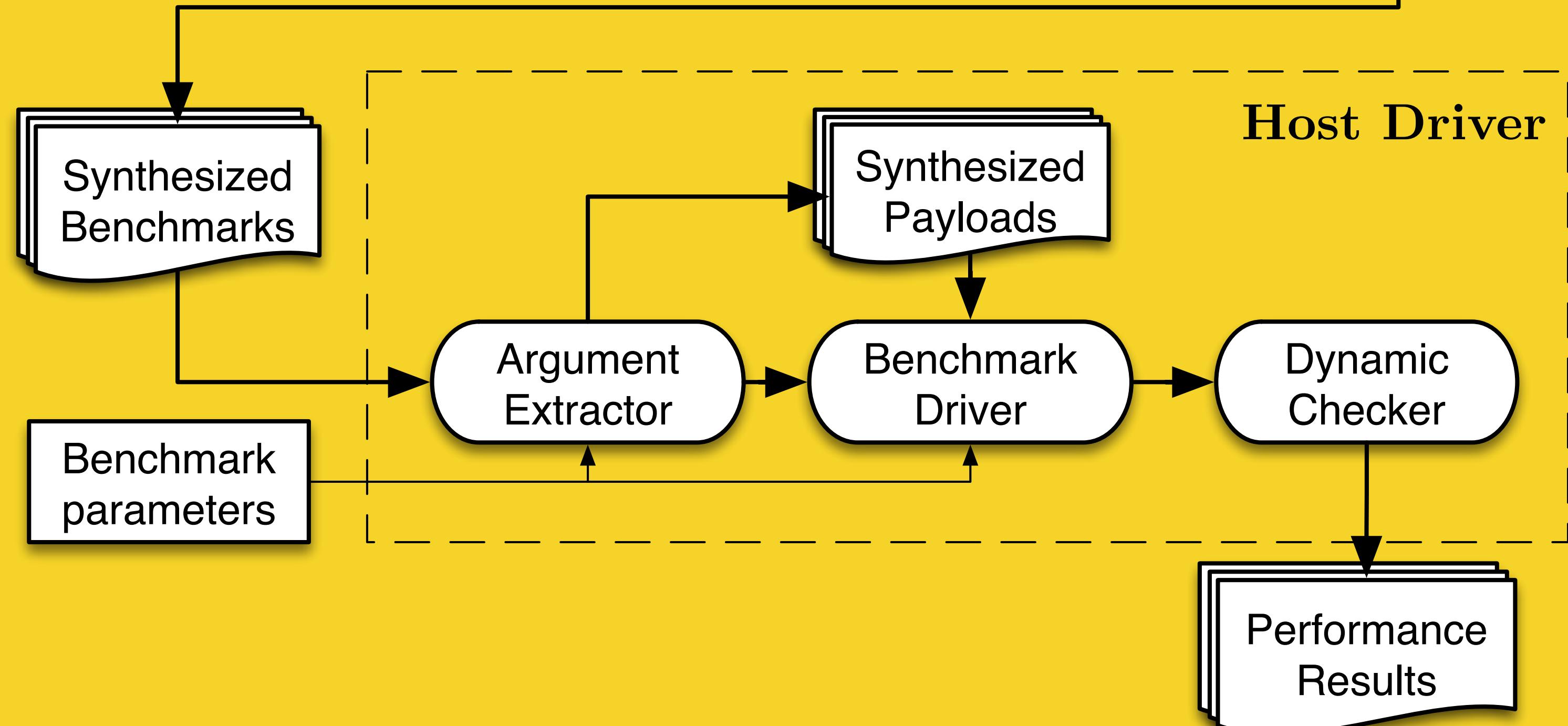
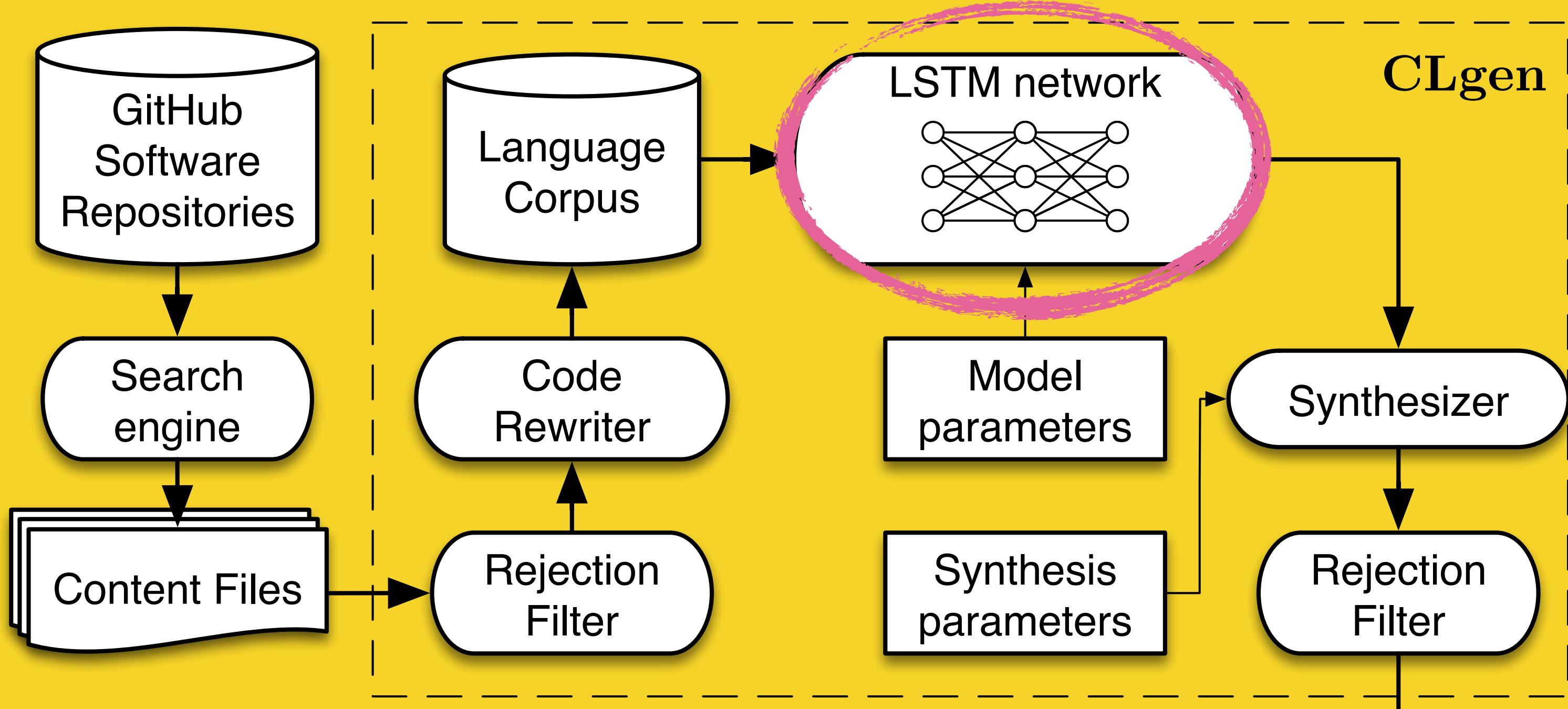
~~Strip comments~~  
~~Preprocess~~  
~~Rewrite function names~~  
~~Rewrite variable names~~  
~~Enforce code style~~

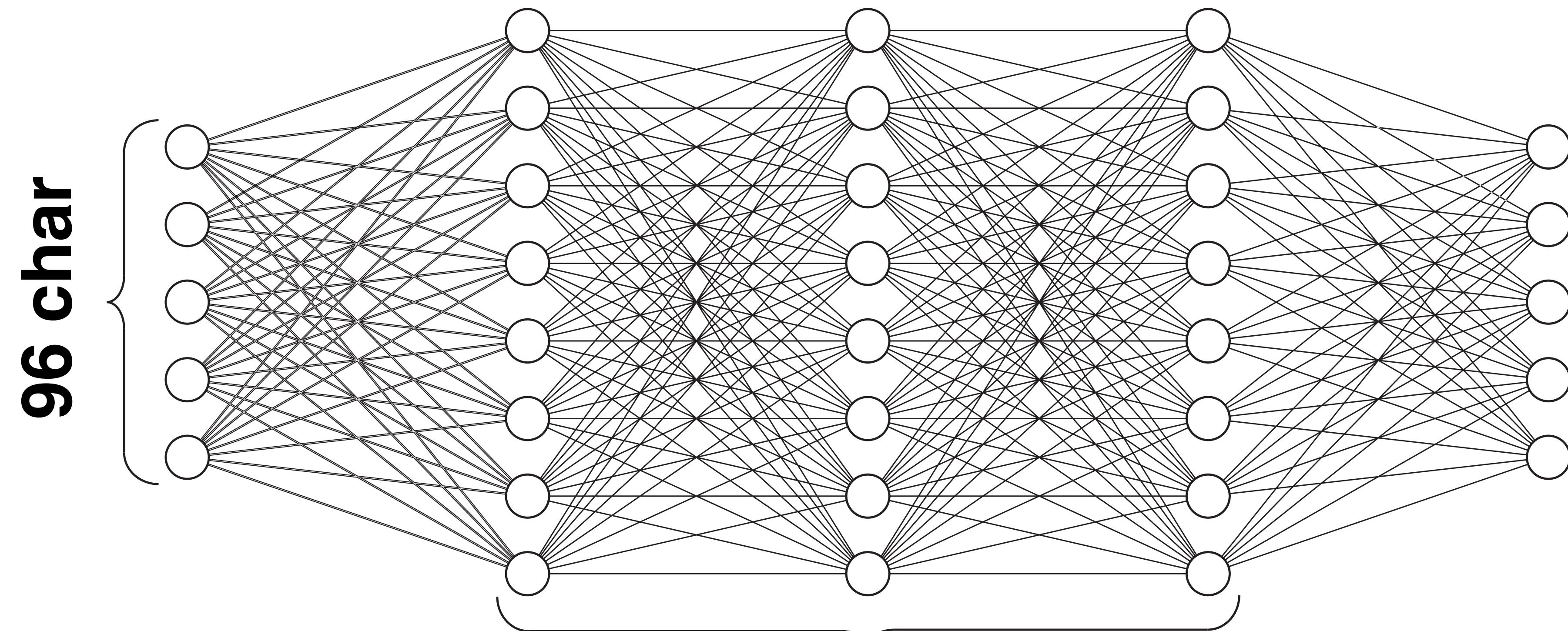
~~Does it compile?~~

~~Does it contain instructions?~~

(33% discard rate)





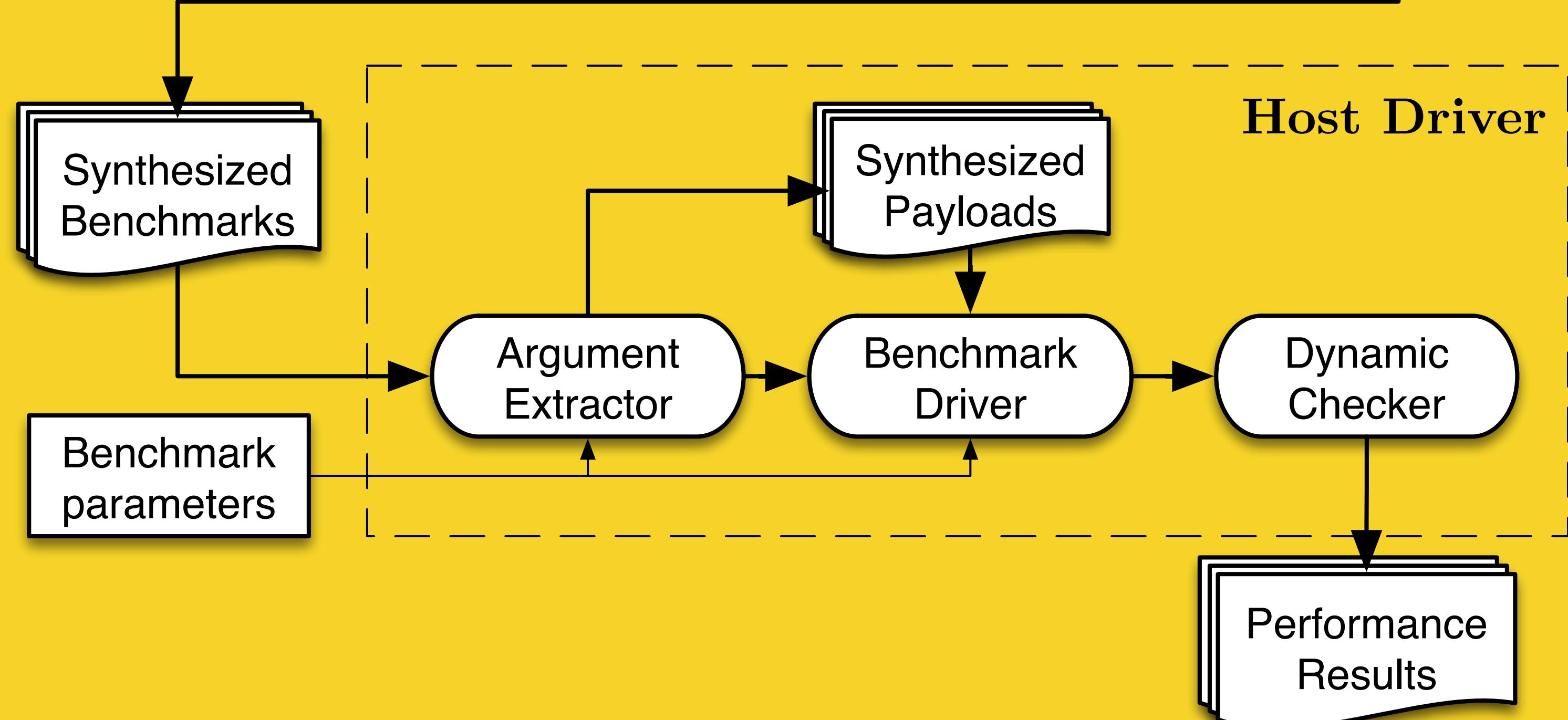
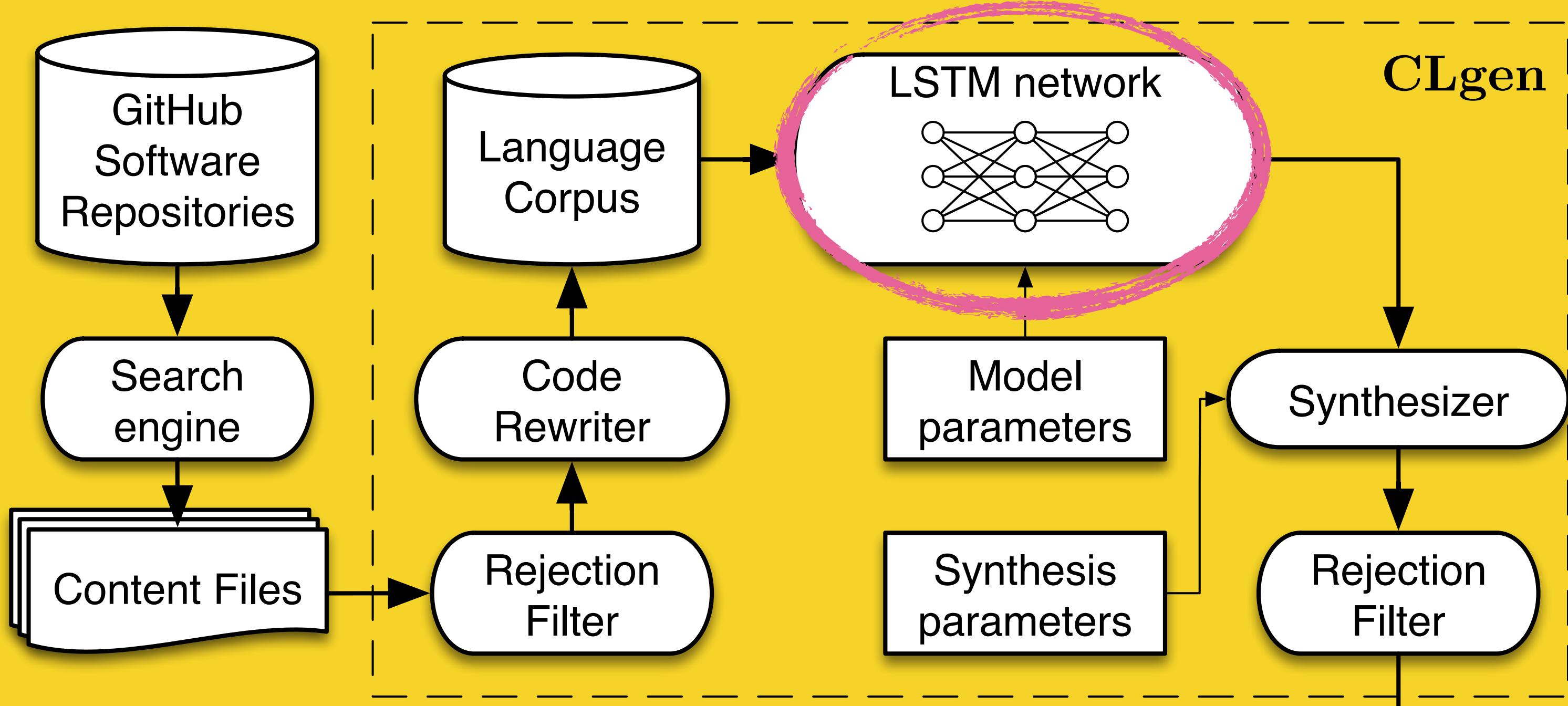


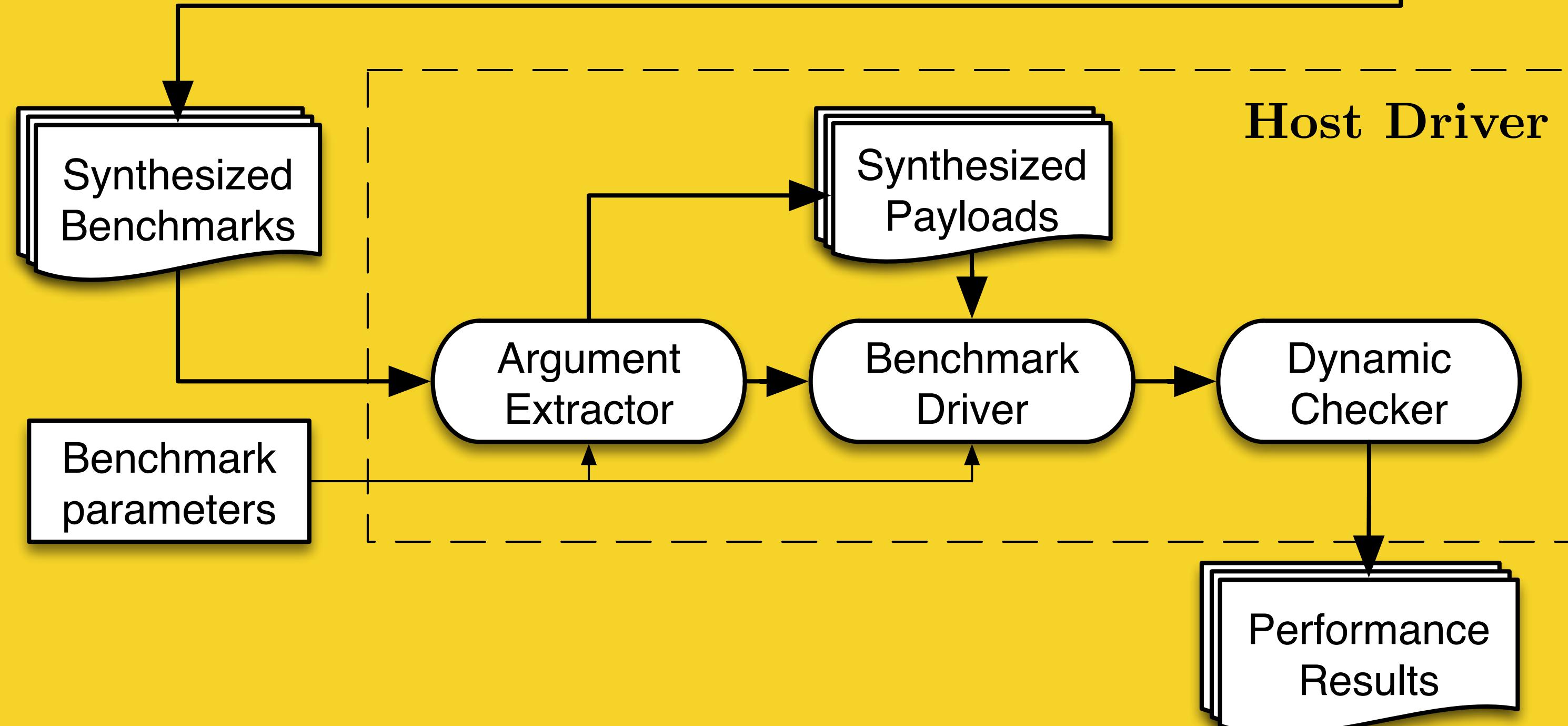
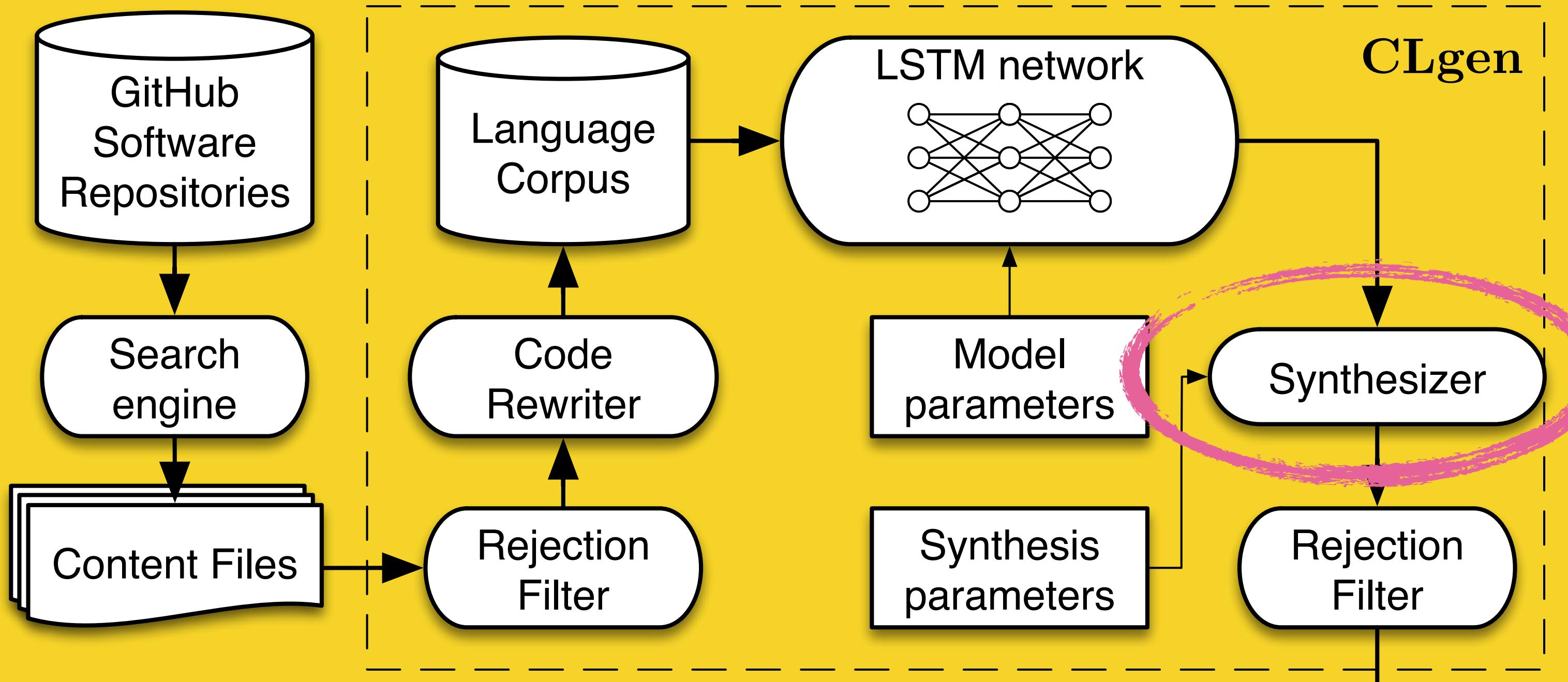
**2048 node, 3 layer LSTM**

## character level language model

$$y = f(x)$$

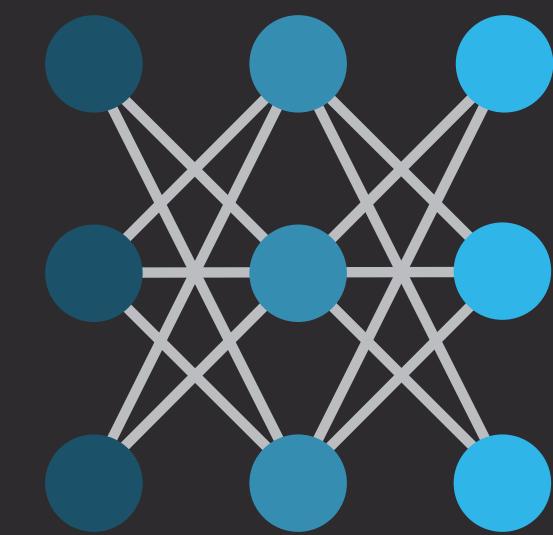
next char      char array





# *kernel synthesis*

```
string S = '__kernel void A(__global float* a) {'  
int depth = 1  
while depth > 0:  
    char c = predict_next_character(S)  
    if c == '{':  
        depth += 1  
    if c == '}':  
        depth -= 1  
    S += c  
  
return S
```



# clgen

Deep Learning Program Generator.

```
$ clgen model.json sampler.json
```

```
{  
    "corpus": {  
        "path": "~/kernels"  
    },  
    "model_type": "lstm",  
    "rnn_size": 2048,  
    "num_layers": 3,  
    "max_epochs": 50  
}
```

```
{  
    "kernels": {  
        "args": [  
            "__global float*",  
            "__global float*",  
            "const int"  
        ]  
    },  
    "sampler": {  
        "max_kernels": 1000  
    }  
}
```

Fork me on GitHub

```
$ clgen model.json sampler.json 2>/dev/null
```

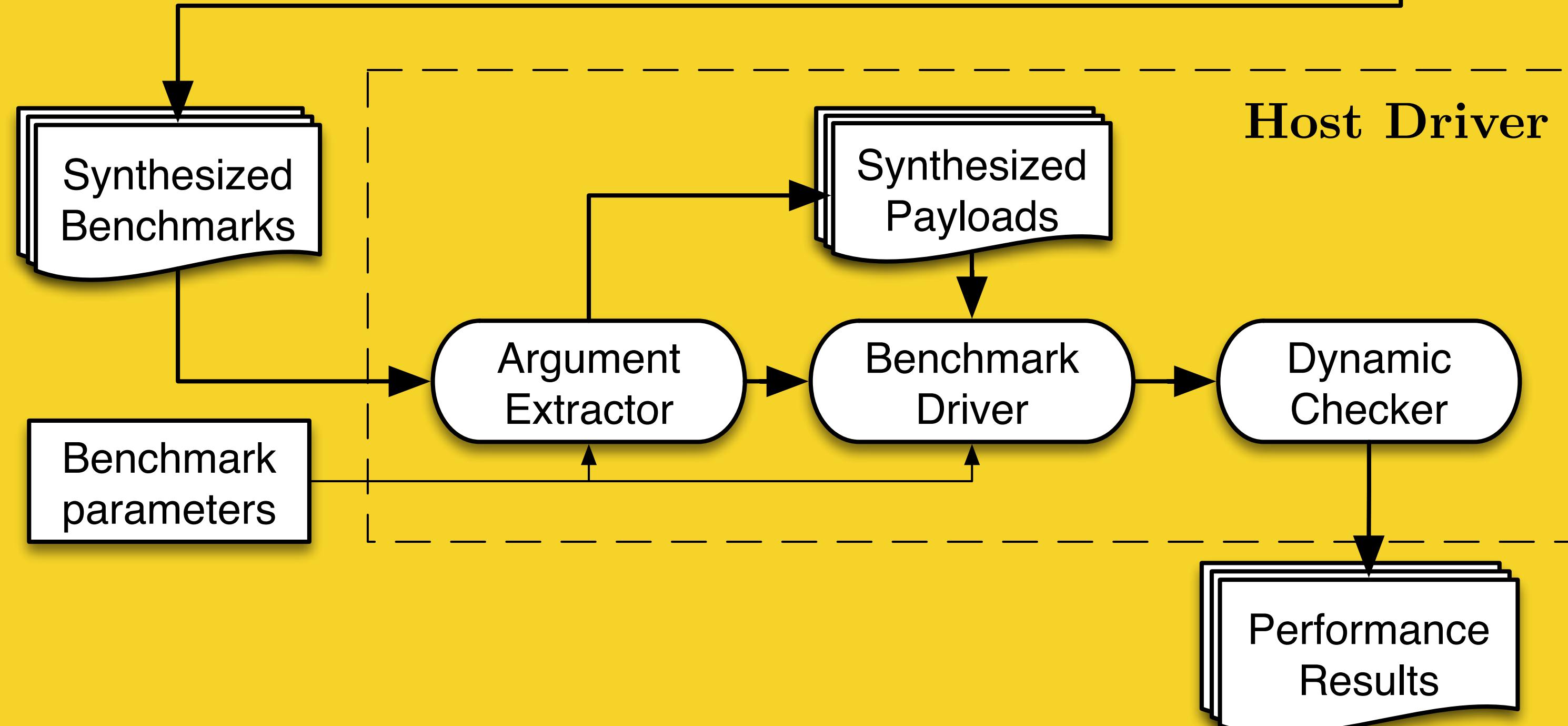
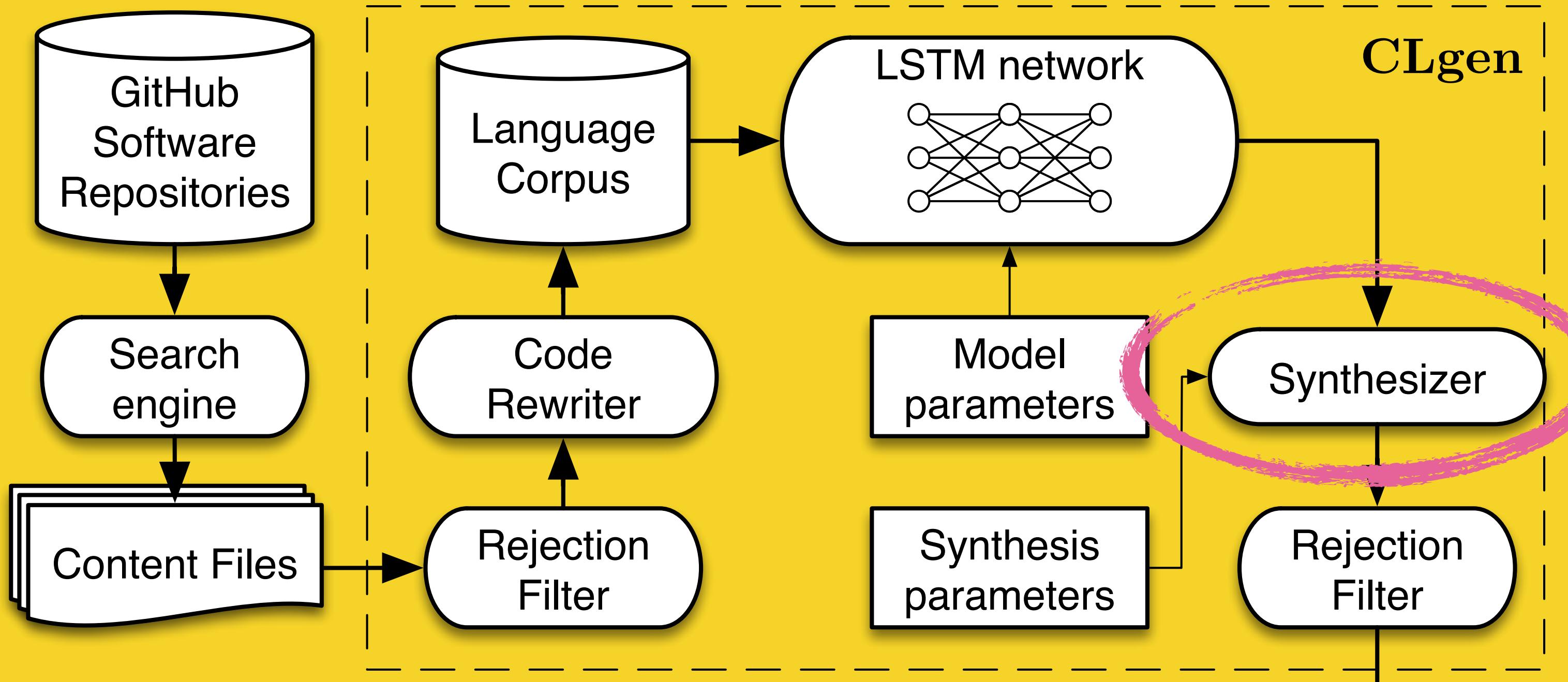
*Fork me on GitHub*

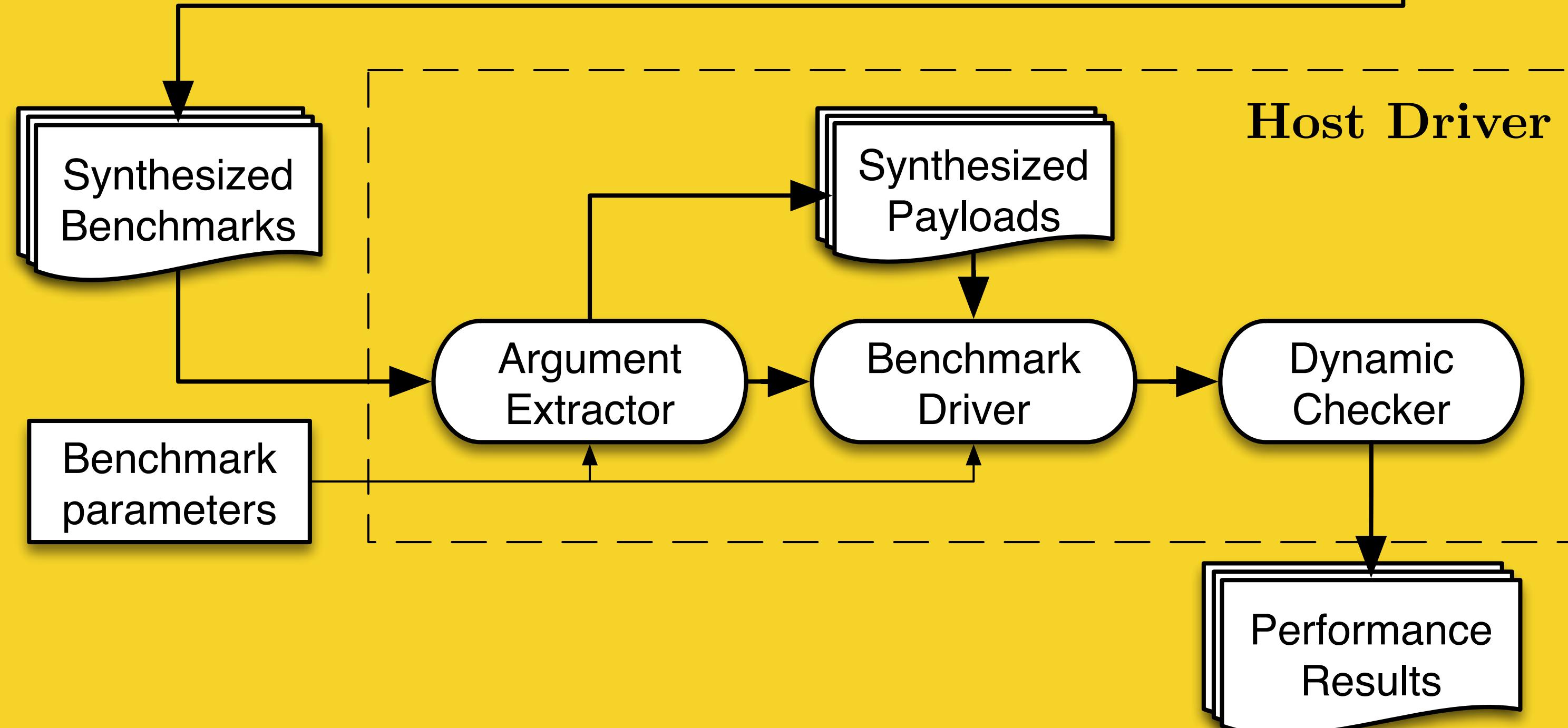
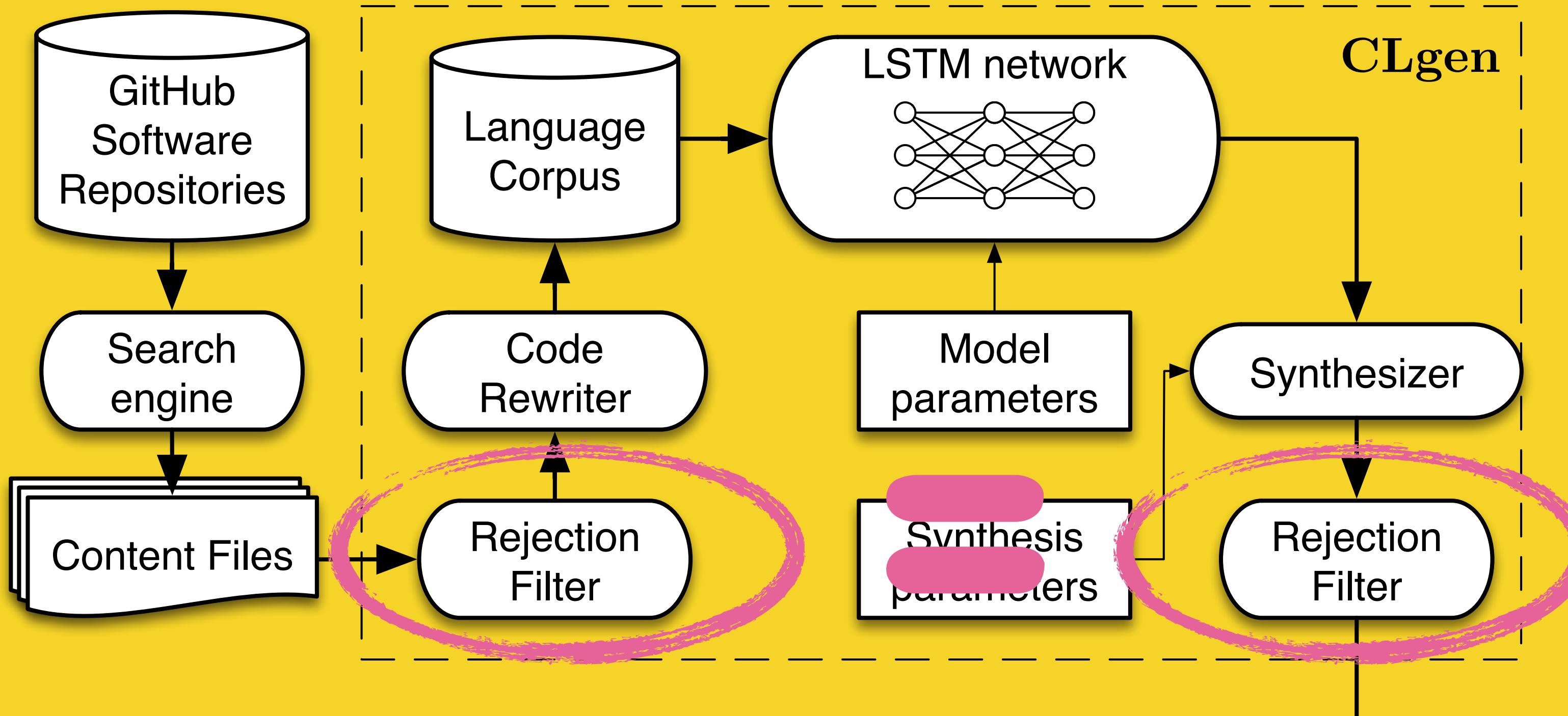
```
__kernel void A(__global float* a,
                __global float* b,
                __global float* c,
                const int d) {
    int e = get_global_id(0);
    float f = 0.0;
    for (int g = 0; g < d; g++) {
        c[g] = 0.0f;
    }
    barrier(1);

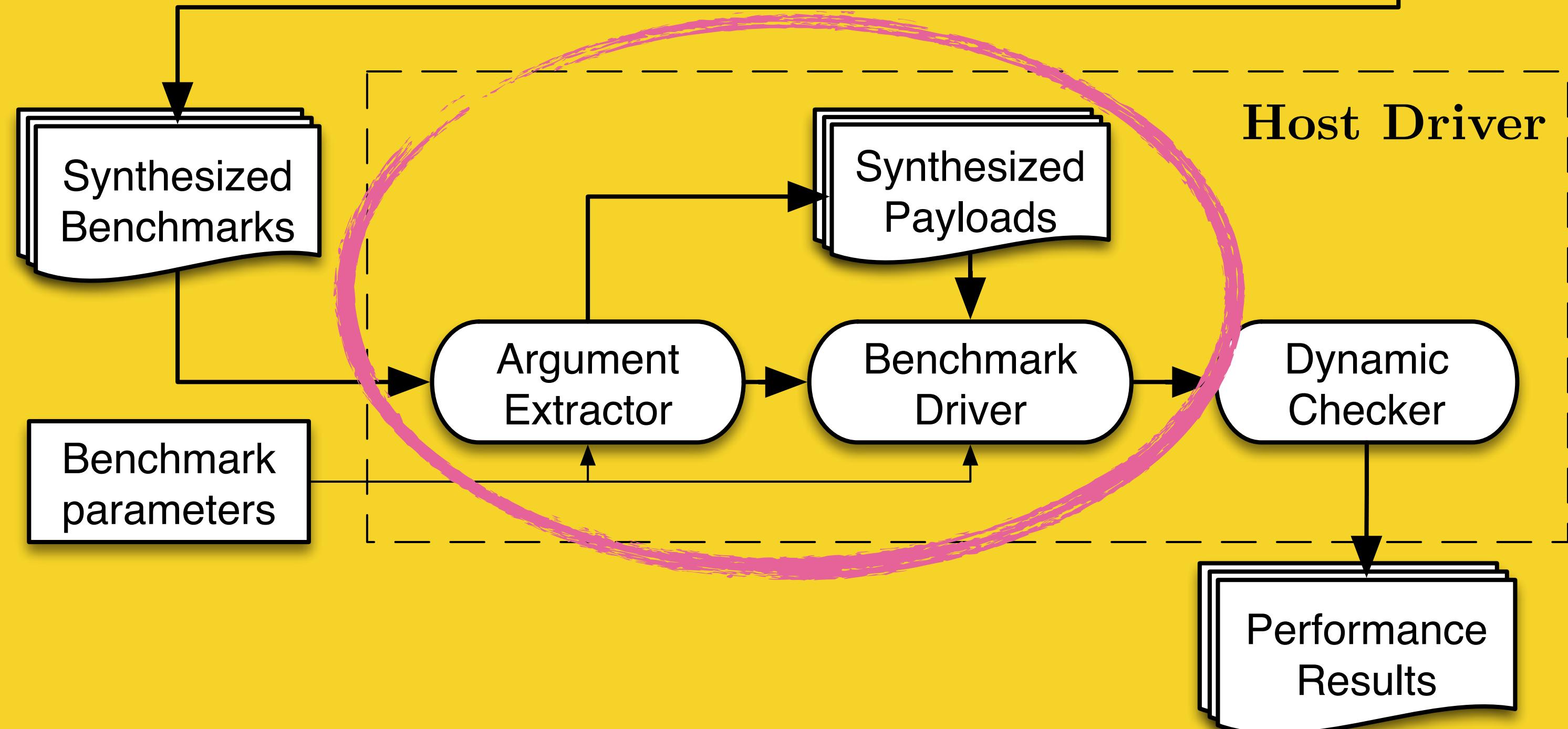
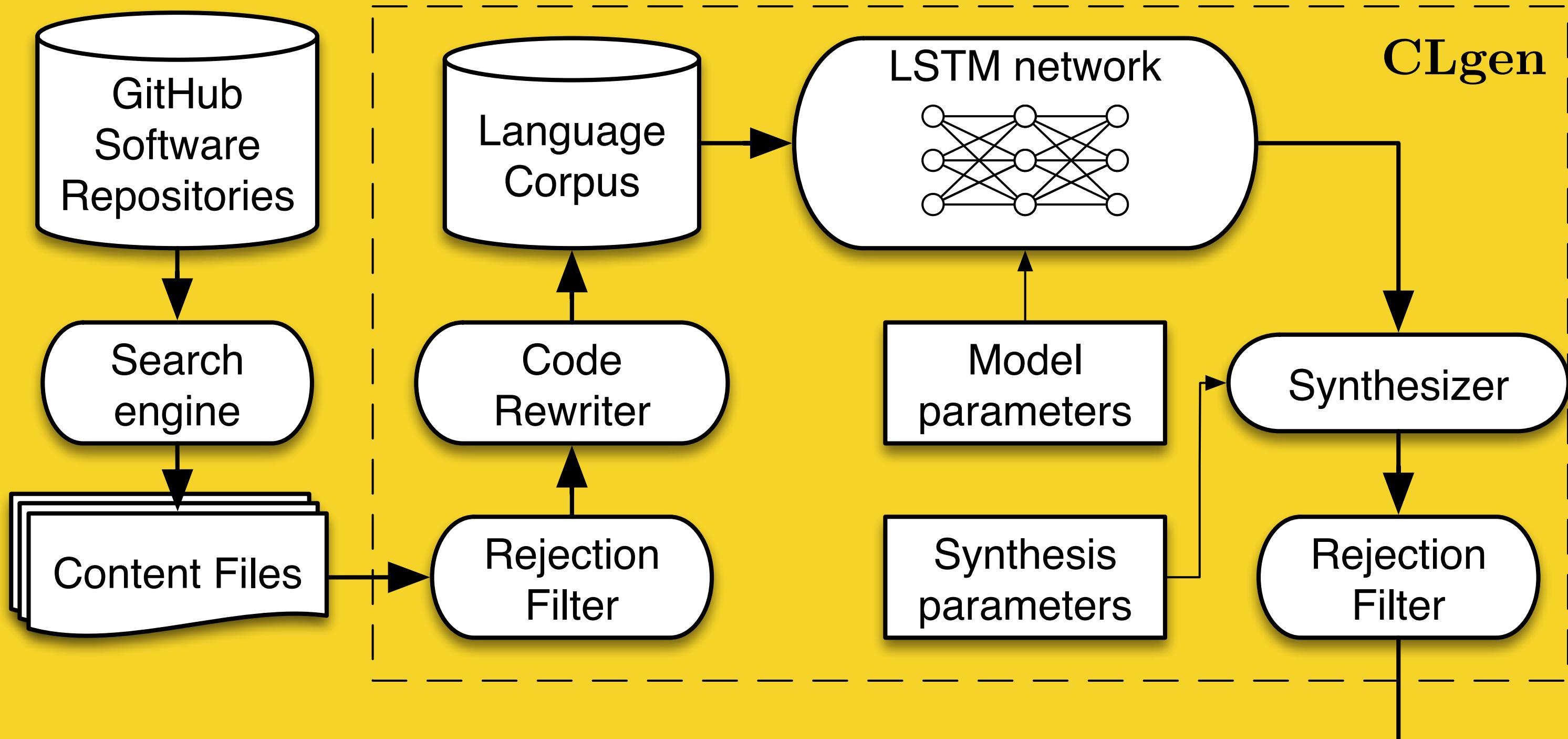
    a[get_global_id(0)] = 2 * b[get_global_id(0)];
}
```

```
__kernel void A(__global float* a,
                __global float* b,
                __global float* c,
                const int d) {
    int e = get_global_id(0);
    if (e >= d) {
        return;
    }
    c[e] = a[e] + b[e] + 2 * a[e] + b[e] + 4;
}
```

```
__kernel void A(__global float* a,
                __global float* b,
                __global float* c,
                const int d) {
    unsigned int e = get_global_id(0);
    float16 f = (float16)(0.0);
    for (unsigned int g = 0; g < d; g++) {
        float16 h = a[g];
        f.s0 += h.s0;
        f.s1 += h.s1;
        /* snip ... */
        f.sE += h.sE;
        f.sF += h.sF;
    }
    b[e] = f.s0 + f.s1 + f.s2 + f.s3 + f.s4 +
           f.s5 + f.s6 + f.s7 + f.s8 + f.s9 + f.sA +
           f.sB + f.sC + f.sD + f.sE + f.sF;
}
```







```
__kernel void A(__global float* a,  
                __global float* b,  
                __global float* c,  
                const float d,  
                const int e) {  
    int f = get_global_id(0);  
    if (f >= e) {  
        return;  
    }  
    c[f] = a[f] + b[f] + 2 * c[f] + d + 4;  
}
```

Payload for size S:

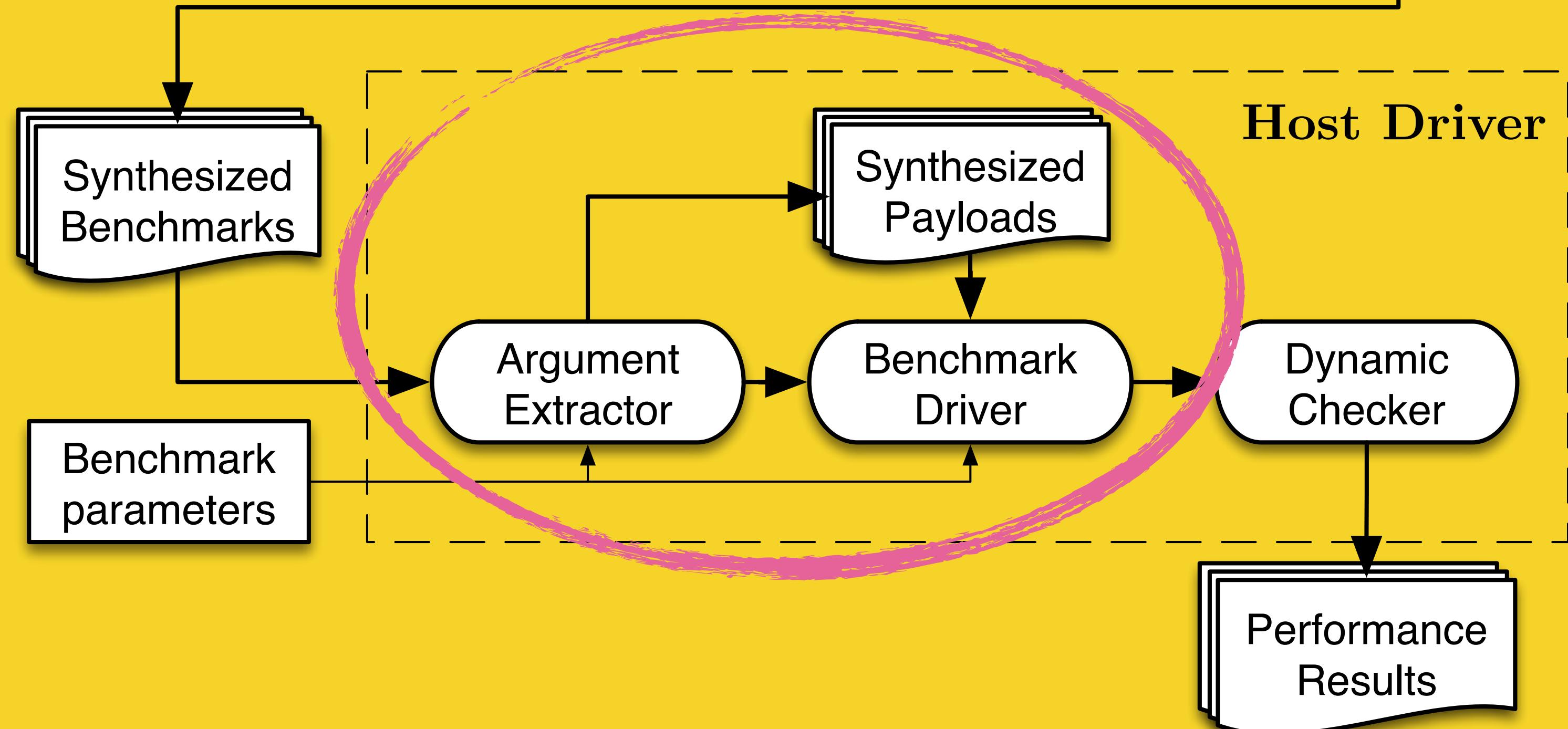
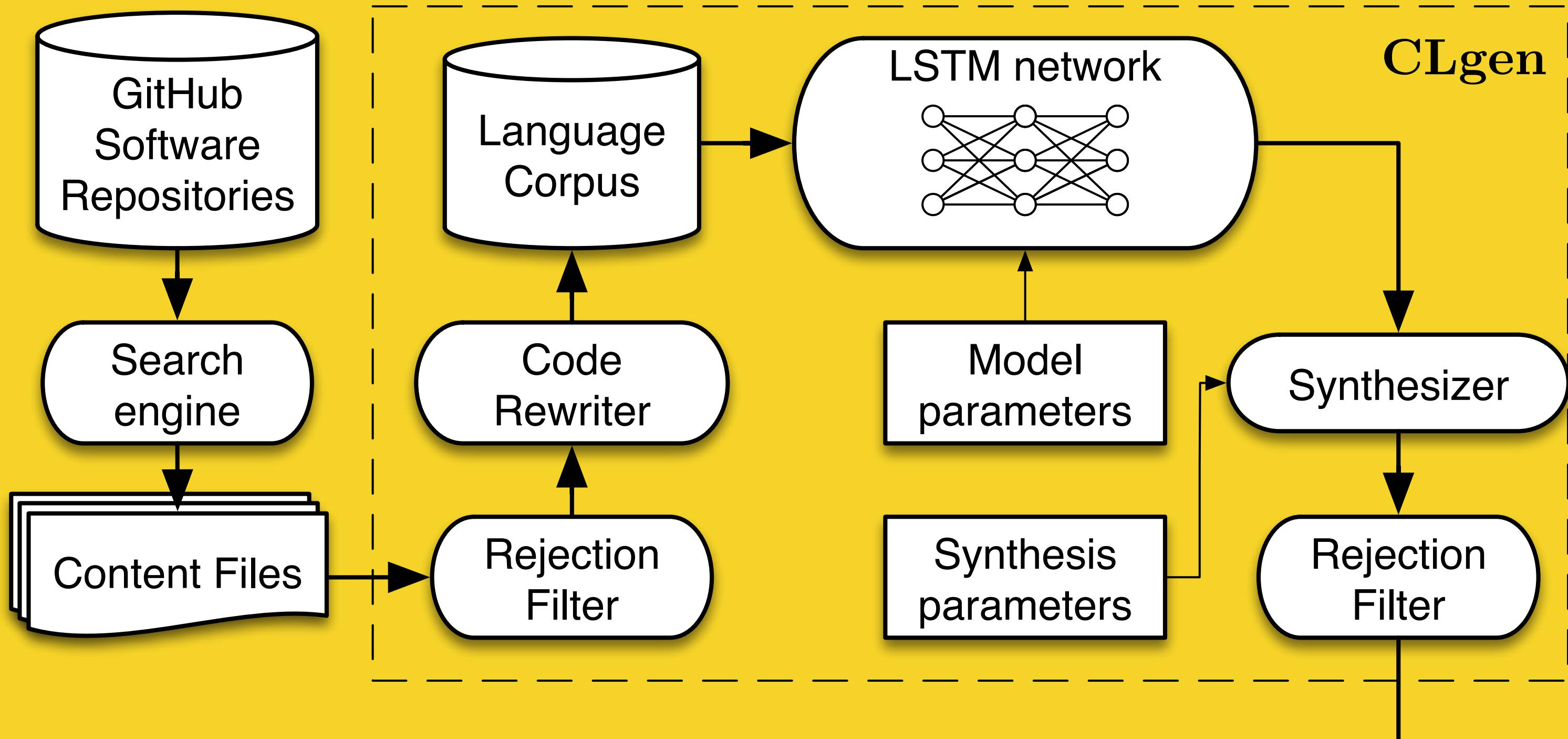
rand() \* [S]

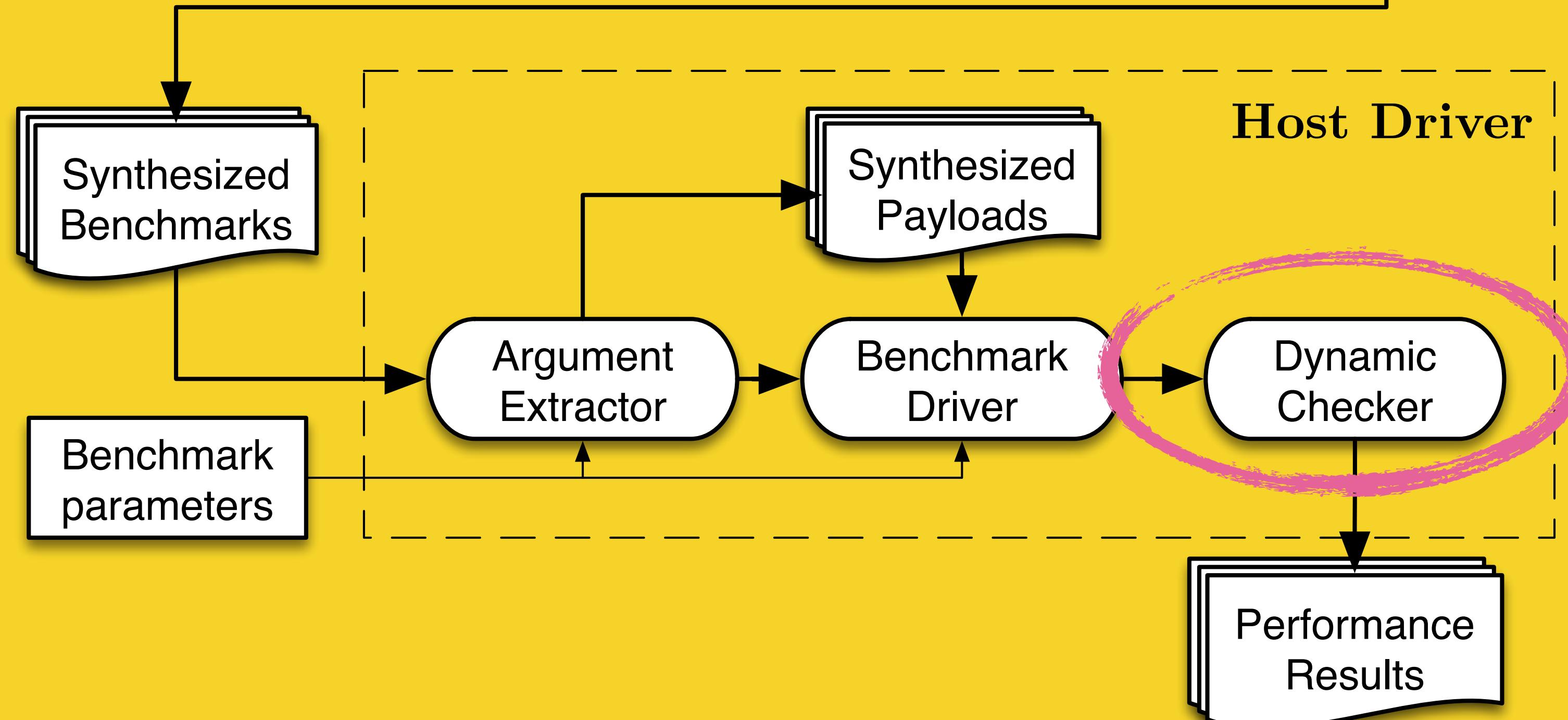
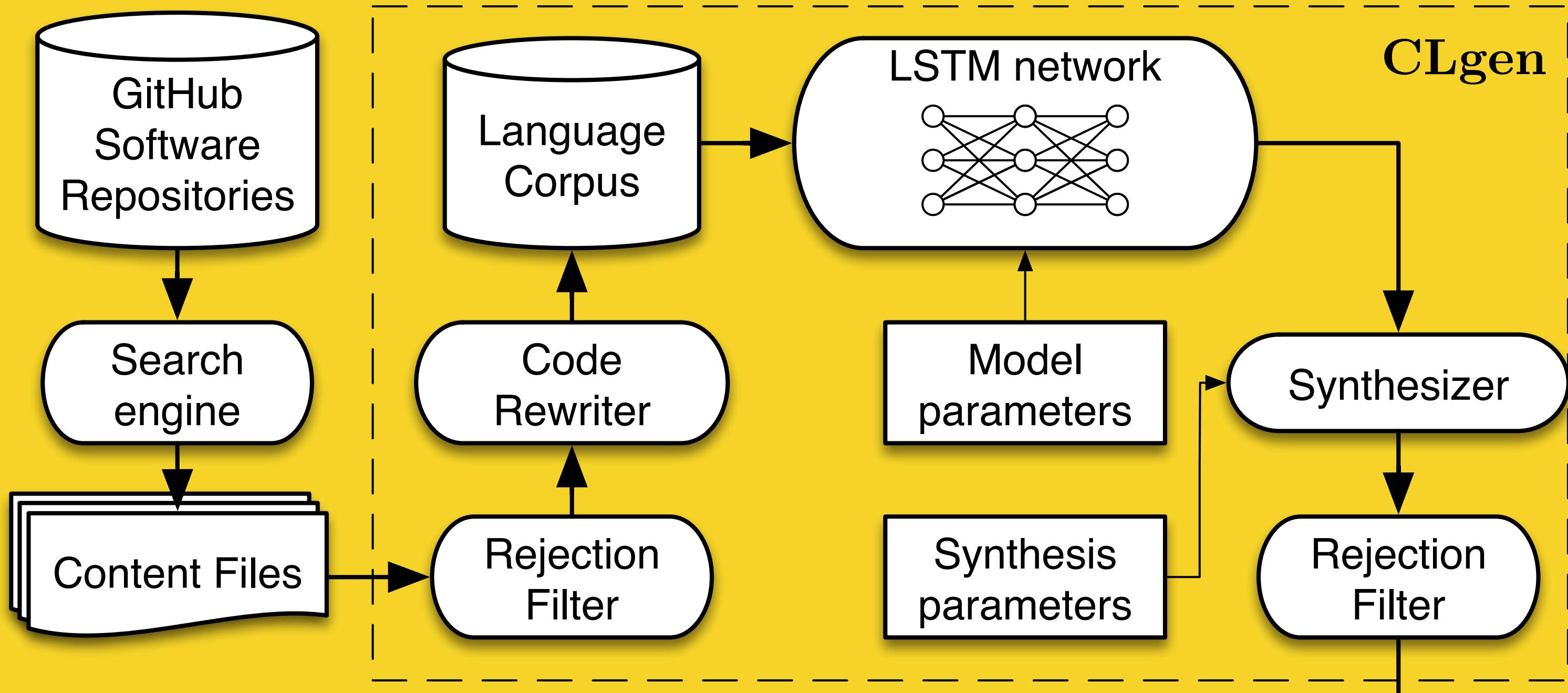
rand() \* [S]

rand() \* [S]

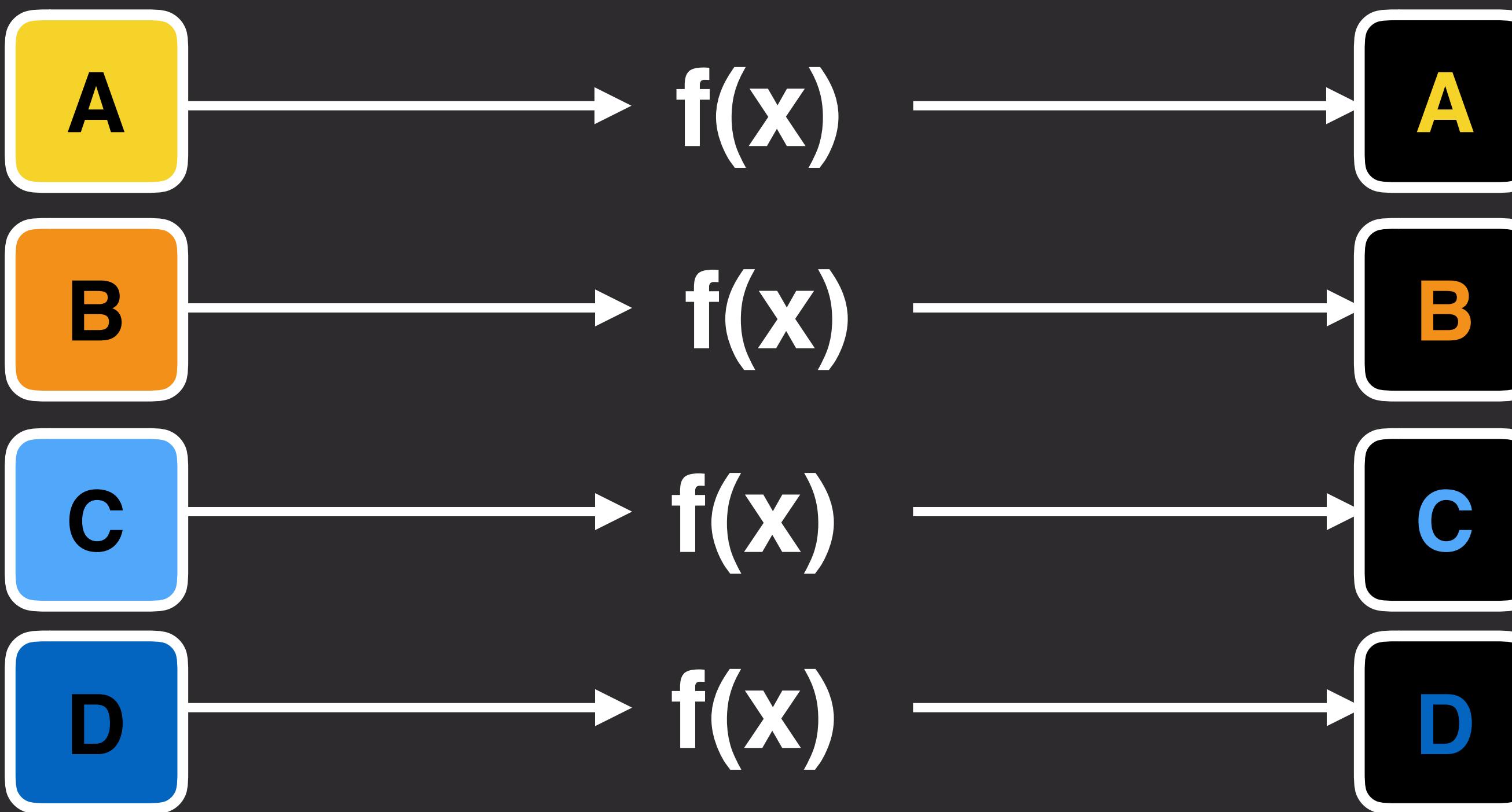
rand()

S



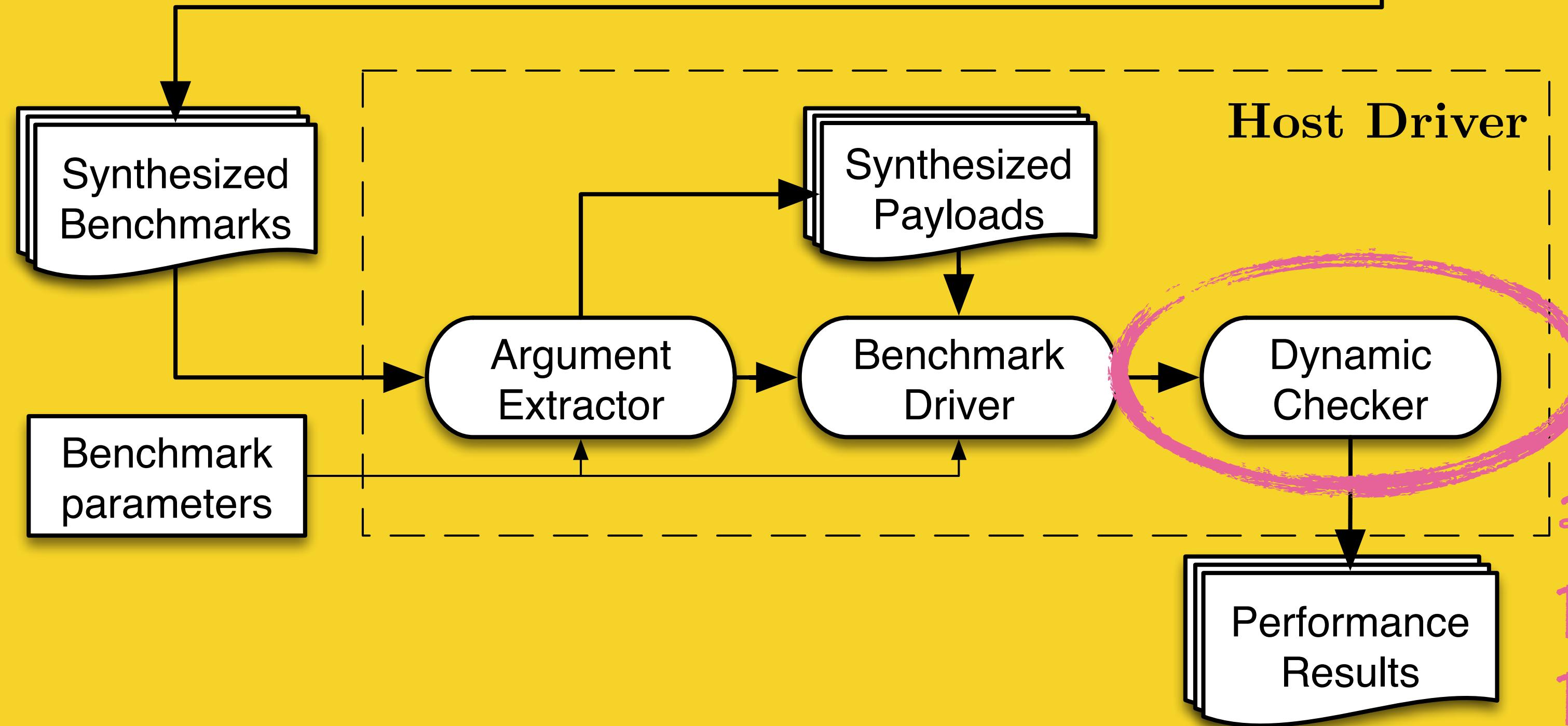
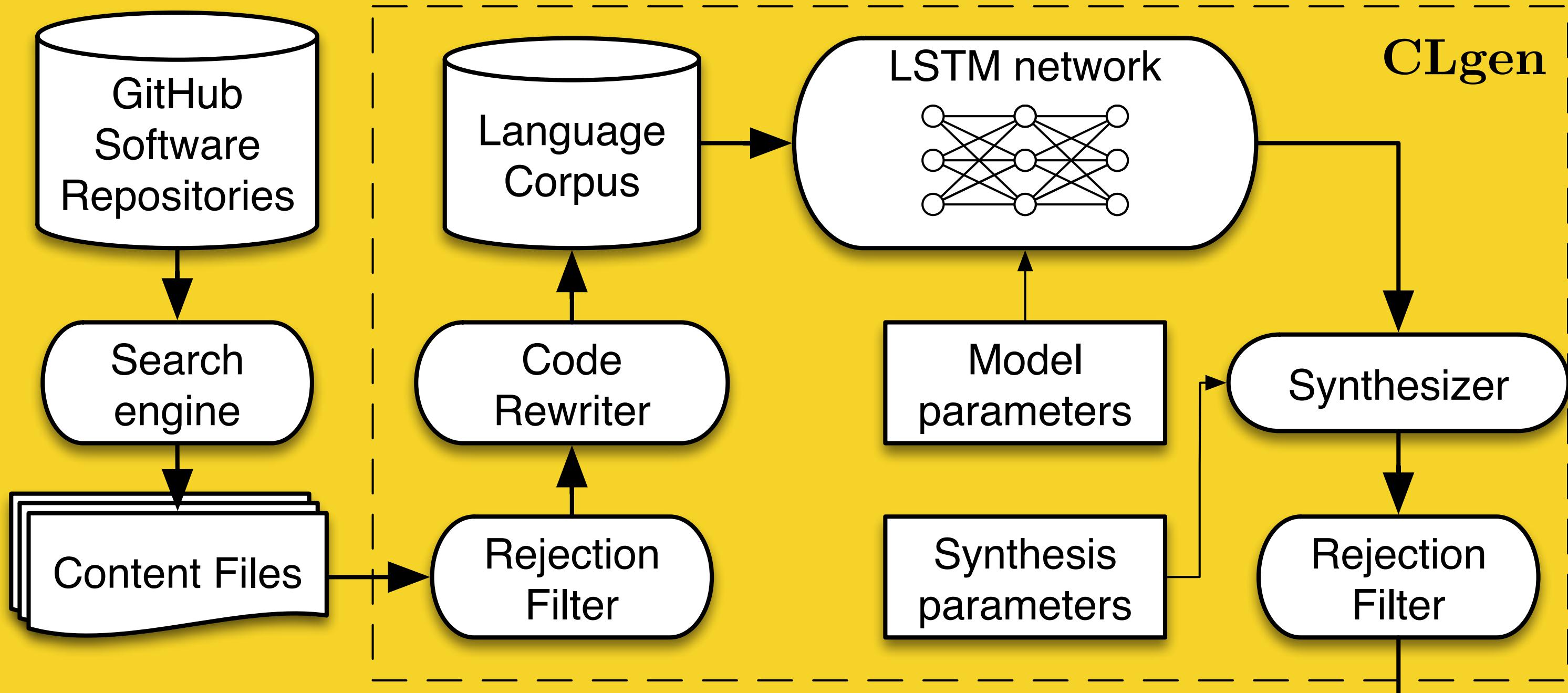


**Generate  
Inputs**



**Verify:**

- $A \neq A$  else *no outputs*
- $A \neq C$  else *input insensitive*
- $A = B$  else *non deterministic*



# results

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) - (0.5f * f + g * 1.0f
9             ↪ / 18.0f + e / 2.0f);
10    }
11    for (c = 0; c < 30; c++) {
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;
```

Listing 7: Sample 7

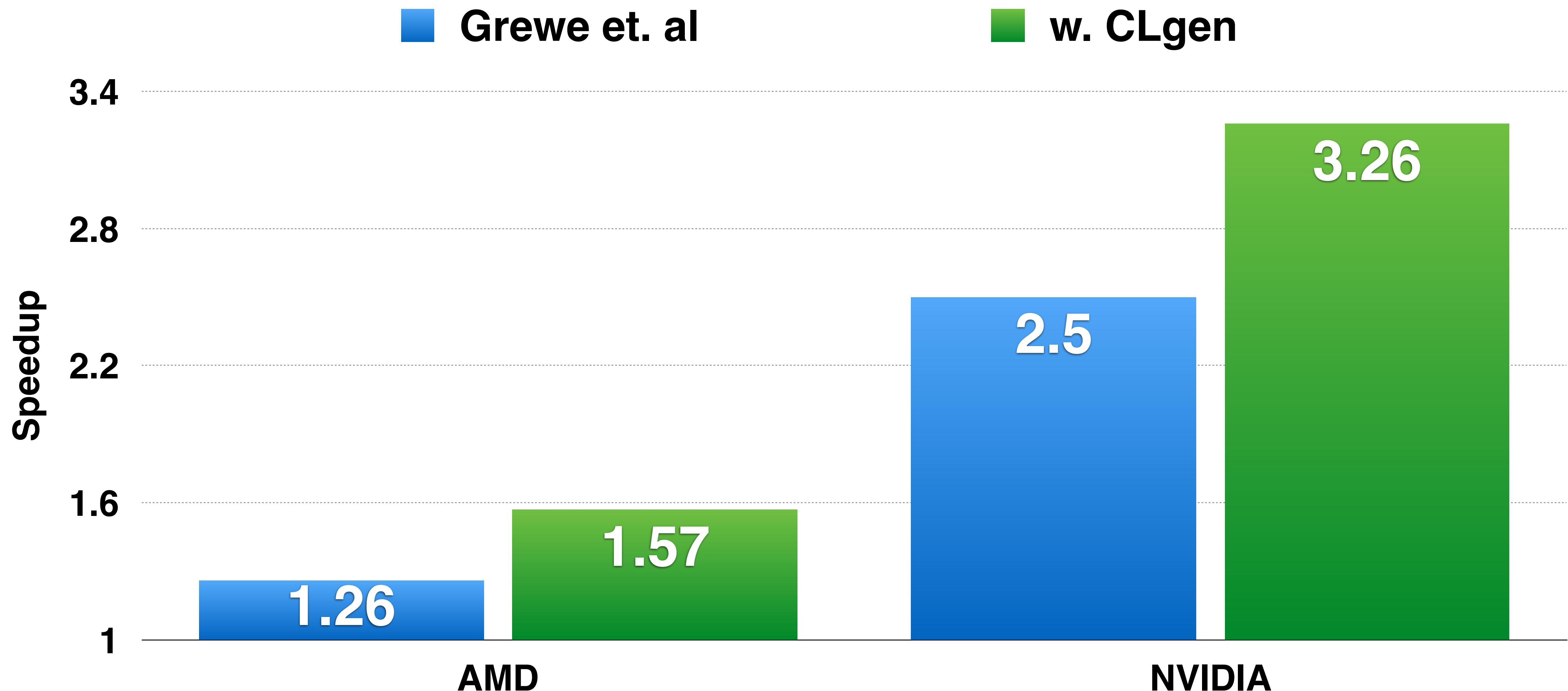
```
1 __kernel void A(int a, int b, int c, __global const float* d, __global
2     ↪ const float* e, __global float* f, float g) {
3     const int h = get_local_id(0);
4     const int i = get_group_id(0);
5
6     const int j = 4 * i + h;
7     const int k = 4 * i + h + a;
8     if (4 * i + h + a < c)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

Listing 10: Sample 10

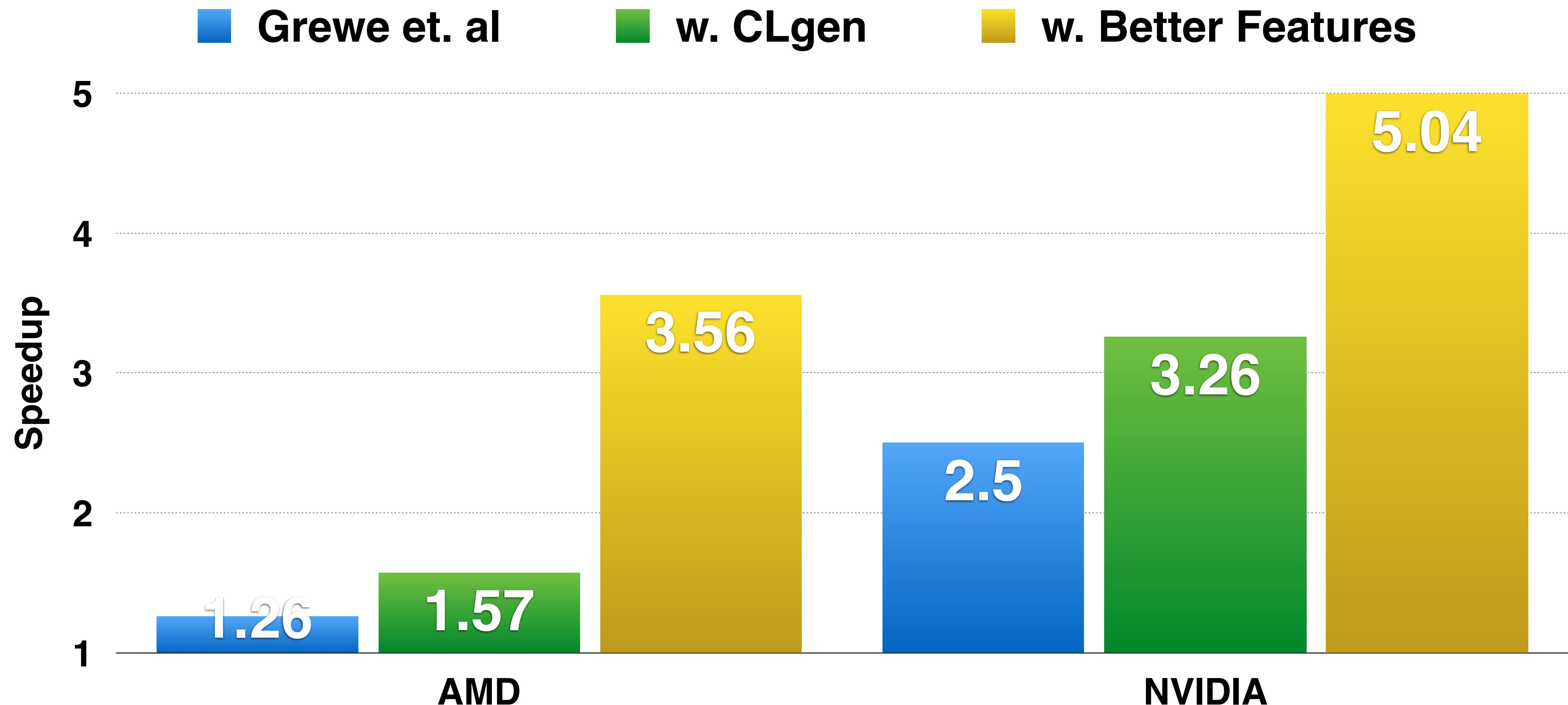
```
1 __kernel void A(__global ulong *a) {
2     int i, j;
3     struct S0 c_8;
4     struct S0* p_7 = &c_8;
5     struct S0 c_9 = {
6         {{0x43250E6DL,2UL},{0x43250E6DL,2UL},{0x43250E6DL,2UL},
7          {0x43250E6DL,2UL},{0x43250E6DL,2UL},{0x43250E6DL,2UL},
8          {0x43250E6DL,2UL},{0x43250E6DL,2UL}},
9         0x4BF90EDCAD2086BDL,
10     };
11     c_8 = c_9;
12     barrier(0 | 1);
13 }
```

52%  
blind test  
<http://humanorrobot.uk>

# 7 benchmarks, 1,000 synthetic benchmarks. 1.27x faster



**71 benchmarks, 1,000 synthetic benchmarks. 4.30x faster**



# *concluding remarks*

**problem: insufficient benchmarks**

**use DL to learn PL semantics and usage**

**turing tested! ;-)**

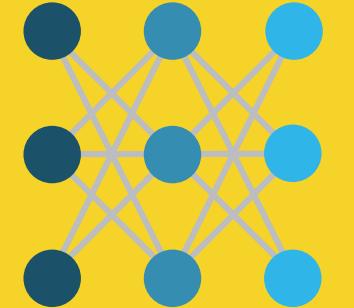
**improved model performance and  
design**

# Synthesizing Benchmarks for Predictive Modeling



For the paper, code and data:  
<http://chriscummins.cc/cgo17>



 **CLgen**  
Deep Learning Program Generator.

For the TensorFlow neural network:  
<http://chriscummins.cc/clgen>

# Extra slides

1. Synthesis algorithm
2. Predictive model features
3. Experimental Setup
4. Model Hyperparameters
5. Limitations & Strengths

---

## Algorithm 1 Sampling a candidate kernel from a seed text.

---

**Require:** LSTM model  $M$ , maximum kernel length  $n$ .

**Ensure:** Completed sample string  $S$ .

```
1:  $S \leftarrow \text{"__kernel void A(const int a) {}}$  Seed text
2:  $d \leftarrow 1$  Initial code block depth
3: for  $i \leftarrow |S|$  to  $n$  do
4:    $c \leftarrow \text{predictcharacter}(M, S)$  Generate new character
5:   if  $c = \{\}$  then
6:      $d \leftarrow d + 1$  Entered code block, increase depth
7:   else if  $c = \}$  then
8:      $d \leftarrow d - 1$  Exited code block, decrease depth
9:   end if
10:   $S \leftarrow S + c$  Append new character
11:  if  $depth = 0$  then
12:    break Exited function block, stop sampling
13:  end if
14: end for
```

---

---

Raw Code Features		
<code>comp</code>	static	#. compute operations
<code>mem</code>	static	#. accesses to global memory
<code>localmem</code>	static	#. accesses to local memory
<code>coalesced</code>	static	#. coalesced memory accesses
<code>transfer</code>	dynamic	size of data transfers
<code>wgsize</code>	dynamic	#. work-items per kernel

(a) Individual code features


---

Combined Code Features	
F1: <code>transfer/(comp+mem)</code>	commun.-computation ratio
F2: <code>coalesced/mem</code>	% coalesced memory accesses
F3: <code>(localmem/mem) × wgsize</code>	ratio local to global mem accesses × #. work-items
F4: <code>comp/mem</code>	computation-mem ratio

(b) Combinations of raw features

---

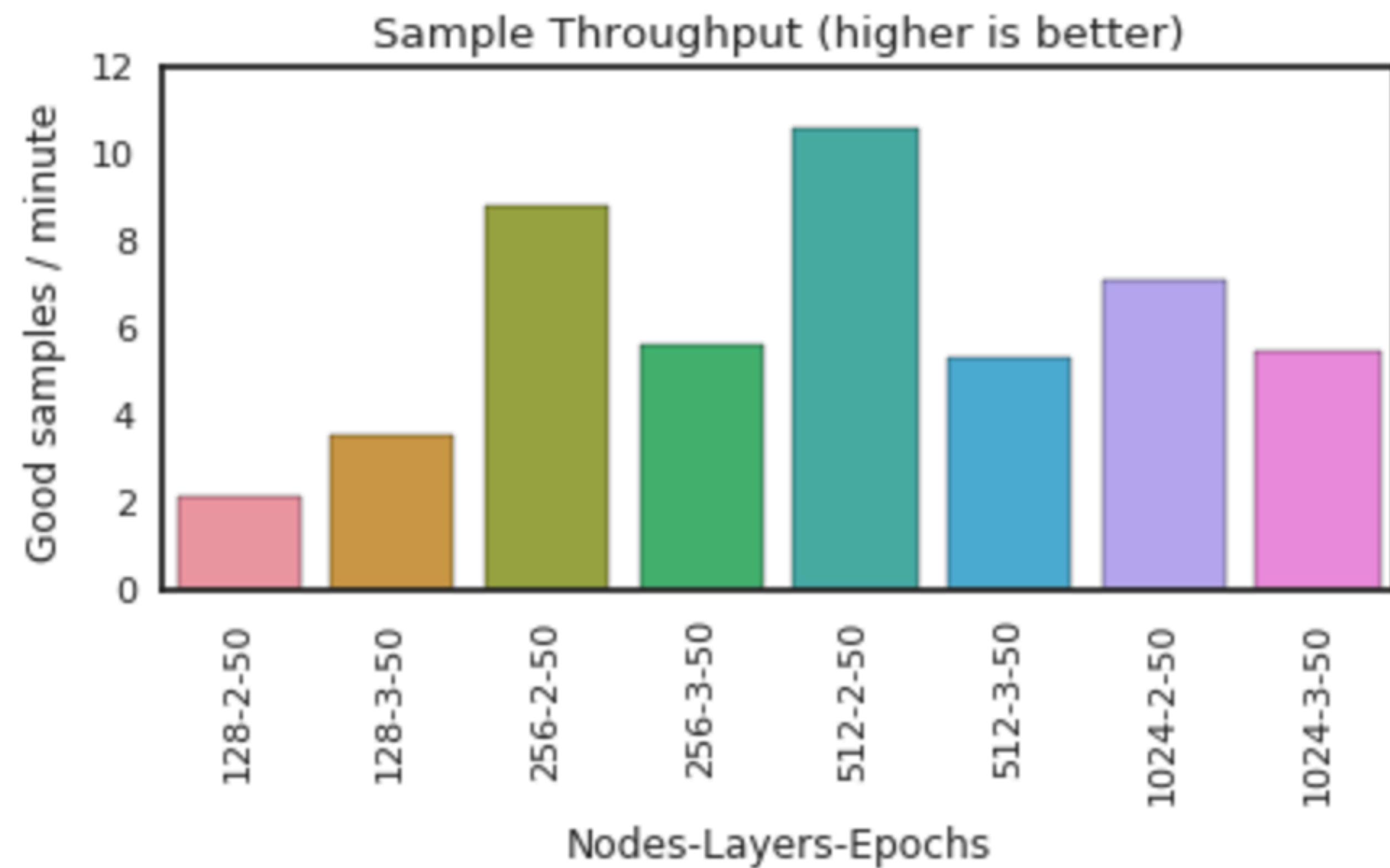
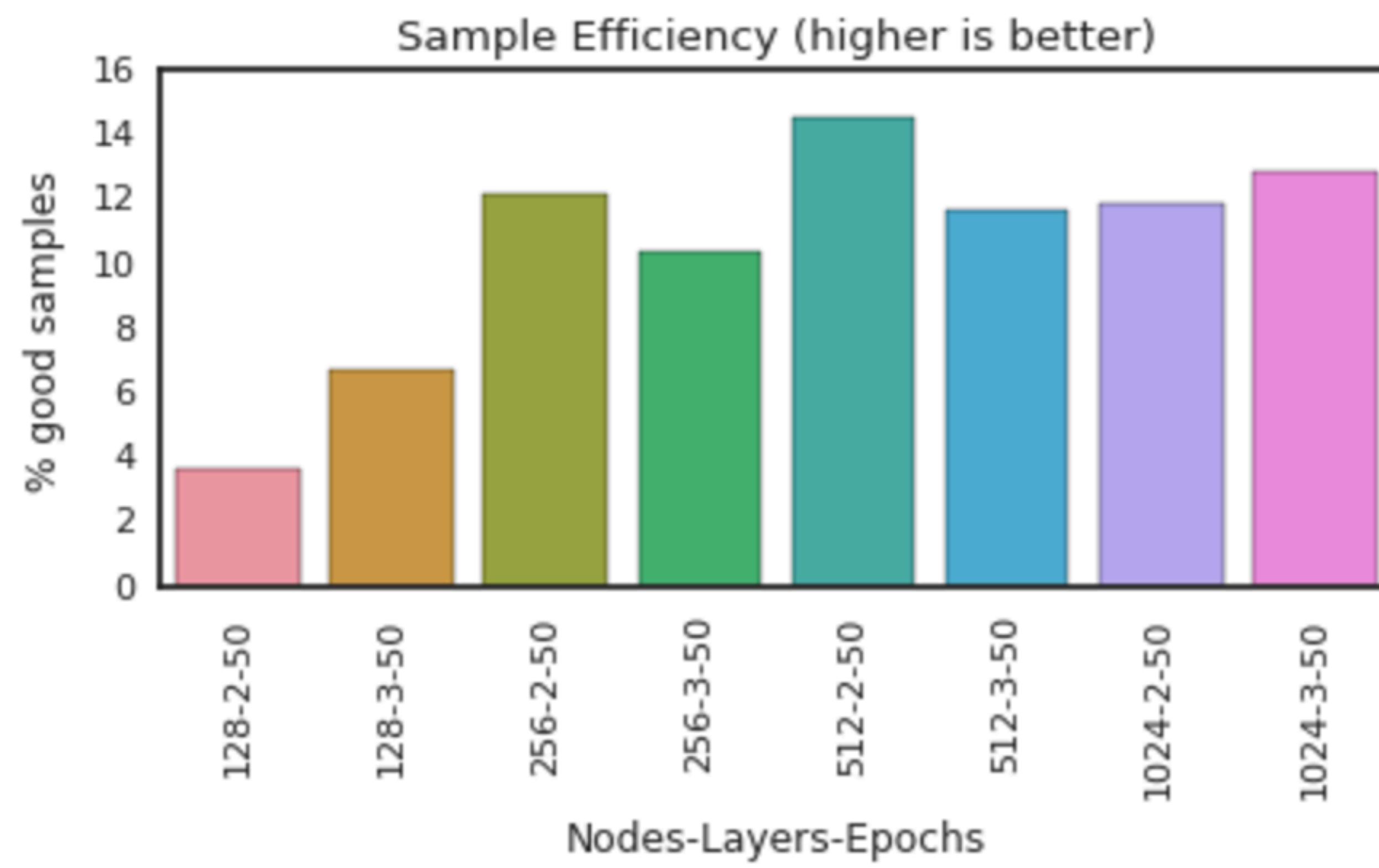
**Table 2.** *Grawe et al.* model features.

	<b>Version</b>	<b>#. benchmarks</b>	<b>#. kernels</b>
<b>NPB (SNU [29])</b>	1.0.3	7	114
<b>Rodinia [30]</b>	3.1	14	31
<b>NVIDIA SDK</b>	4.2	6	12
<b>AMD SDK</b>	3.0	12	16
<b>Parboil [31]</b>	0.2	6	8
<b>PolyBench [32]</b>	1.0	14	27
<b>SHOC [33]</b>	1.1.5	12	48
<b>Total</b>	-	71	256

**Table 3.** List of benchmarks.

	<b>Intel CPU</b>	<b>AMD GPU</b>	<b>NVIDIA GPU</b>
<b>Model</b>	Core i7-3820	Tahiti 7970	GTX 970
<b>Frequency</b>	3.6 GHz	1000 MHz	1050 MHz
<b>#. Cores</b>	4	2048	1664
<b>Memory</b>	8 GB	3 GB	4 GB
<b>Throughput</b>	105 GFLOPS	3.79 TFLOPS	3.90 TFLOPS
<b>Driver</b>	AMD 1526.3	AMD 1526.3	NVIDIA 361.42
<b>Compiler</b>	GCC 4.7.2	GCC 4.7.2	GCC 5.4.0

**Table 4.** Experimental platforms.



# Limitations

Rejection sampling is wasteful.  
No support for things declared outside of kernel scope.

# Strengths

Basically\* language agnostic.  
35 million repos on GitHub. We're using 0.00004%.  
Generates 2000 OpenCL benchmarks per machine per day.