## **Glenside** Partitioning Deep Learning Hardware and Software

Gus Smith, University of Washington PL/Arch External Talk Series @ CAPRA, May 27th, 2020

#### Who?

- 2nd year at UW working in PL/ Arch with Zach Tatlock/Luis Ceze
- Interested in how we can advance architecture w/ PL techniques





## Architects accelerate deep learning with custom hardware.







## Architects accelerate deep learning with custom hardware.



CADE METZ 04.05.17 12:03 PM





Building an Al Chip Saved Google From Building a Dozen New Data Centers





## ...but designing deep learning stacks is complex!

#### **In-Datacenter Performance Analysis of a Tensor Processing Unit<sup>TM</sup>**

Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon Google, Inc., Mountain View, CA USA Email: {jouppi, cliffy, nishantpatil, davidpatterson} @google.com

To appear at the 44th International Symposium on Computer Architecture (ISCA), Toronto, Canada, June 26, 2017.

#### TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems (Preliminary White Paper, November 9, 2015)

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng Google Research\*



W

Workloads





Workloads











#### (hardware engineers)



(software engineers)



#### An Oversimplified View of HW–SW stack design Workloads We need matrix multiplication! We can build it! And we'll (hardware engineers) (software engineers) build ReLU, why not?









#### An Oversimplified View of HW–SW stack design Workloads We need matrix multiplication! We can build it! And we'll (hardware engineers) (software engineers) build ReLU, why not? ...we'll figure out Hardware how to use it! Intrinsics





#### An Oversimplified View of HW–SW stack design Workloads We need matrix multiplication! We can build it! And we'll (hardware engineers) (software engineers) build ReLU, why not? ...we'll figure out Hardware how to use it! Compiler Passes Intrinsics







(software engineers)

Compiler Passes





Compiler Passes

Kernels



This system makes sense because chips and compilers are massive projects.





(software engineers)

Compiler Passes

Kernels

This system makes sense because chips and compilers are massive projects.

...but we can't help but notice the inefficiencies!



(software engineers)

Compiler Passes

Kernels

This system makes sense because chips and compilers are massive projects.

...but we can't help but notice the inefficiencies!

• Slow



(software engineers) Compiler Passes

Kernels

This system makes sense because chips and compilers are massive projects.

...but we can't help but notice the inefficiencies!

- Slow
- Misses design points



(software engineers) Compiler Passes Kernels

This system makes sense because chips and compilers are massive projects.

...but we can't help but notice the inefficiencies!

- Slow
- Misses design points
- Design knowledge is lost in communication!







## Why not design the hardware and software at the same time?

### 



Workloads









### 





Workloads

(automatic codesigner)



l see a 16x16 vector–

matrix multiplication...

Workloads

(automatic codesigner)



l see a 16x16 vector–

matrix multiplication...





l see a 16x16 vector–

matrix multiplication...

l see a 40x102 fused

vector-matrix multiplication/

ReLU...





l see a 16x16 vector–

matrix multiplication...

l see a 40x102 fused

vector-matrix multiplication/

ReLU...





l see a 16x16 vector–

matrix multiplication...

l see a 40x102 fused

vector-matrix multiplication/

ReLU...

Hardware Intrinsics



#### Kernels

Compiler Passes



#### ...this is not a new idea. **Hardware–software codesign** is a dream that people have had for some time.

have had for some time.

Hardware–software codesign often raises a healthy amount of skepticism due to its meager success in the past.

#### ...this is not a new idea. Hardware-software codesign is a dream that people

have had for some time.

Hardware–software codesign often raises a healthy amount of skepticism due to its meager success in the past.

One of the primary issues is that the space of potential hardware-software designs is massive!

#### ...this is not a new idea. Hardware-software codesign is a dream that people

### 



Workloads






Workloads

(automatic codesigner)



l see a 16x16 vector–

matrix multiplication...

Workloads

(automatic codesigner)



l see a 16x16 vector–

matrix multiplication...





I see a 16x16 vector– matrix multiplication...

...or 4x4...





I see a 16x16 vector– matrix multiplication...

...or 4x4...

































## If only we had a tool for representing massive search spaces!







Remy Wang Chandrakana Nandi **Oliver** Flatt Zach Tatlock Pavel Panchekha

### Using E-Graphs for Optimization with egg

000 N

#### Max Willsey github.com/mwillsey/egg





#### (a\*2)/2

x*2	-	x << 1
((x*y)*z)	-	(x *(y * z))
x/x	-	l
x*1	-	×
x*y	-	y*x











### (a\*2)/2and (a >> 1)/2are represented a





eclass sizes S1, S2, S3...

∏s<sub>i</sub> terms Es<sub>i</sub> space





 $x^*y \leftrightarrow y^*x$ 









## =



### Uses E-Graphs to simplify expressions Total speed up: 2.77x

Simplify speed up: 5–60x

Better results!

HERBIE







### =

#### Based on Apache SystemML

#### Replace the optimization of LA expressions with E-Graph

#### Extraction aware of density/memory-usage

#### 1.2x - 5x better results









Remy Wang Chandrakana Nandi **Oliver** Flatt Zach Tatlock Pavel Panchekha

### Using E-Graphs for Optimization with egg

000 N

#### Max Willsey github.com/mwillsey/egg

### perfect for egraphs!

Fundamentally, this is a problem of generating many equivalent programs —

Fundamentally, this is a problem of generating **many equivalent programs**— perfect for egraphs!

However, in our case, "programs" will simultaneously represent both hardware *and* software!

- Fundamentally, this is a problem of generating **many equivalent programs** perfect for egraphs!
- However, in our case, "programs" will simultaneously represent both hardware *and* software!
- Rewrites over our programs then become a natural way to represent changing the partition between hardware and software!

Glenside's goal: For a given deep learning workload, efficiently explore the massive space of accelerator designs using egraphs.

# Glenside Design

#### of hardware atoms

**Given:** (1) A deep learning model in a high level deep learning DSL; (2) a library

**Given:** (1) A deep learning model in a high level deep learning DSL; (2) a library of *hardware atoms* 

1. Convert model to our egraph representation

of hardware atoms

- 1. Convert model to our egraph representation
- 2. Run *rewrites* until saturation or timeout

**Given:** (1) A deep learning model in a high level deep learning DSL; (2) a library
of hardware atoms

- 1. Convert model to our egraph representation
- 2. Run *rewrites* until saturation or timeout
  - Software  $\rightarrow$  software

of hardware atoms

- 1. Convert model to our egraph representation
- 2. Run *rewrites* until saturation or timeout
  - Software  $\rightarrow$  software
  - Software  $\rightarrow$  hardware

of hardware atoms

- 1. Convert model to our egraph representation
- 2. Run *rewrites* until saturation or timeout
  - Software  $\rightarrow$  software
  - Software  $\rightarrow$  hardware
- 3. Extract hardware–software program from egraph

of hardware atoms

- 1. Convert model to our egraph representation
- 2. Run *rewrites* until saturation or timeout
  - Software  $\rightarrow$  software
  - Software  $\rightarrow$  hardware
- 3. Extract hardware-software program from egraph 4. Separate hardware-software program into hardware description and software
- schedule

**Given:** (1) A deep learning model in a high level deep learning DSL; (2) a library of hardware atoms

- 1. Convert model to our egraph representation
- 2. Run *rewrites* until saturation or timeout
  - Software  $\rightarrow$  software
  - Software  $\rightarrow$  hardware
- 3. Extract hardware-software program from egraph
- schedule

4. Separate hardware-software program into hardware description and software

































# How do we represent vector-matrix multiplication in an egraph?



The power of egraphs comes from their grouping of equivalent program nodes (enodes) into the same equality class (eclass).

(enodes) into the same equality class (eclass).

another (referential transparency).

- The power of egraphs comes from their grouping of equivalent program nodes
- For this to be possible, we need a representation where **every node has a value** and there are no side-effects, so that equivalent nodes can be swapped for one

(vec-mat-mult a b) := (map-dot-product (cartesian-product a (cols b))

(map-dot-product
 (cartesian-product
 a
 (cols b)
)



(map-dot-product

#### (cartesian-product

а (cols b)





••••

#### (cartesian-product

#### (map-dot-product

(cartesian-product (map-do a (cols b) )



1x32

### Hardware Atoms

### Hardware Atoms

Hardware atoms are small hand-designed units which we compose to build a hardware design.

### Hardware Atoms

Hardware atoms are small hand-designed units which we compose to build a hardware design.

Each hardware atom is paired with a functional description which we can use to map it in to the workload.

A systolic array is one way to implement vector – matrix multiplication in hardware.

A systolic array is one way to implement vector-matrix multiplication in hardware.

They have a fixed size input they can handle, e.g. 16x16 (or 256x256 on the TPU!)

A systolic array is one way to implement vector – matrix multiplication in hardware.

TPU!)

We describe its functional behavior with a **rewrite**:

They have a fixed size input they can handle, e.g. 16x16 (or 256x256 on the

A systolic array is one way to implement vector – matrix multiplication in hardware.

- TPU!)
- We describe its functional behavior with a **rewrite**:

(map-dot-product

They have a fixed size input they can handle, e.g. 16x16 (or 256x256 on the

- A systolic array is one way to implement vector-matrix multiplication in hardware.
- They have a fixed size input they can handle, e.g. 16x16 (or 256x256 on the TPU!)
- We describe its functional behavior with a rewrite:
- (map-dot-product
  - (cartesian-product x{1xN vector} (cols y{NxM tensor}) ))

- A systolic array is one way to implement vector-matrix multiplication in hardware.
- They have a fixed size input they can handle, e.g. 16x16 (or 256x256 on the TPU!)
- We describe its functional behavior with a rewrite:
- (map-dot-product
  - (cartesian-product x{1xN vector} (cols y{NxM tensor}) ))
  - => (systolic-array N M x y)

# Codesigning a Vector– Matrix Multiply

#### **Goal:** explore hardware – software implementations of

## **Goal:** explore hardware-software implementations of (map-dot-product

### **Goal:** explore hardware – software implementations of (map-dot-product (cartesian-product a{1x32}

# Goal: explore hardware-software implementations of (map-dot-product (cartesian-product a{1x32}

a{1x32} (cols b{32x32})))

# Goal: explore hardware-software implementations of (map-dot-product (cartesian-product a{1x32}
Goal: explore hardware-(map-dot-product (cartesian-product

#### Goal: explore hardware-software implementations of

(map-dot-product (cartesian-product a{1x32})

Using the atom library composed of: (systolic-array N M x y) Where N,M = 16 or 32.

#### **Goal:** explore hardware – software implementations of

(cols b{32x32})))



We have: (map-dot-product (cartesian-product a{1x32}

uct a{1x32} (cols b{32x32})))





### can introduce our hardware atom:

To explore hardware–software splits, we want to find places where we



#### can introduce our hardware atom: (map-dot-product

To explore hardware–software splits, we want to find places where we



To explore hardware—software splits, we want to find places where we can introduce our hardware atom: (map-dot-product (cartesian-product x{1xN vector} (cols y{NxM tensor})))

- => (systolic-array N M x y)
- can introduce our hardware atom:

(map-dot-product

1x32 а



To explore hardware–software splits, we want to find places where we

(cartesian-product x{1xN vector} (cols y{NxM tensor}) ))

- => (systolic-array N M x y)
- (map-dot-product
- can introduce our hardware atom:



To explore hardware–software splits, we want to find places where we

(cartesian-product x{1xN vector} (cols y{NxM tensor}) ))



- => (systolic-array N M x y)
- (map-dot-product
- can introduce our hardware atom:















# We use rewrites to expose places where hardware can be mapped in!











+







+

















a => (concat (slice a ...) (slice a ...))

























#### (let's keep this graph simple!)



#### Recall what we're looking for:






























We won't find them here...the concats are "hiding" the tensor slices!

Can we move the concats somehow?











egg::rewrite!("bubble-concat-through-cols"; "(cols (concat a b))"













### From this, we can extract a hardware description and software schedule. (our next step!)



## Early Results

```
let program = "
 (map-dot-product
  (cartesian-product
   v-32
   (cols t-32-32)
 11
.parse()
.unwrap();
```

### Early Results

#### **let** program = " (map-dot-product (cartesian-product v-32 (cols t-32-32) .parse() .unwrap();

t	r	W	s	=	=	V	e	С	!	[				
	r	e١	NI	ri	i t	e	s	:	•	s	р	l	i	t
	r	e١	NI	ri	i t	e	s	:	•	s	р	l	i	t
	r	e١	NI	ri	i t	e	s	:	•	С	0	l	l	а
	r	e١	NI	ri	i t	e	s	:	•	b	u	b	b	l
	r	e١	NI	ri	i t	e	s	:	•	b	u	b	b	ι
	r	e١	NI	ri	i t	e	s	:	•	b	u	b	b	ι
	r	e١	NI	ri	i t	:e	s	:	•	b	u	b	b	ι
	r	e١	NI	ri	i t	e	s	:	•	b	u	b	b	ι
	r	e١	NI	ri	i t	e	s	:	:	b	u	b	b	ι
	r	e١	NI	ri	i t	:e	s	:	:	b	u	b	b	l
	r	e١	NI	ri	i t	:e	s	:	:	b	u	b	b	l
	r	e١	NI	ri	i t	:e	s	:	•	b	u	b	b	l
	r	e١	NI	ri	i t	e	s	•	•	s	у	s	t	0

### Early Results

```
(0, 16, true),
(1, 16, true),
pse_nested_slices(),
e_concat_through_rows_axis_0(),
e_concat_through_rows_axis_1(),
e_concat_through_cols_axis_0(),
e_concat_through_cols_axis_1(),
e_concat_through_cartesian_product_not_last_axis_left(),
e_concat_through_cartesian_product_not_last_axis_right(),
e_concat_through_cartesian_product_last_axis(),
e_concat_through_map_dot_product_not_last_axis(),
e_concat_through_map_dot_product_last_axis(),
lic_array_vector_matrix(),
```



#### let program = " (map-dot-product (cartesian-product v-32 (cols t-32-32) .parse() .unwrap();

et	rws = vec![
	rewrites::split
	rewrites::split
	rewrites::colla
	rewrites::bubbl
	rewrites::systo

let (egraph, id) = egg::EGraph::<Language, Meta>::from\_expr(&program); **let** runner = egg::Runner::new().with\_egraph(egraph).run(&rws);

### Early Results

(0, 16, **true**), (1, 16, **true**), pse\_nested\_slices(), e\_concat\_through\_rows\_axis\_0(), e\_concat\_through\_rows\_axis\_1(), e\_concat\_through\_cols\_axis\_0(), e\_concat\_through\_cols\_axis\_1(), e\_concat\_through\_cartesian\_product\_not\_last\_axis\_left(), e\_concat\_through\_cartesian\_product\_not\_last\_axis\_right(), e\_concat\_through\_cartesian\_product\_last\_axis(), e\_concat\_through\_map\_dot\_product\_not\_last\_axis(), e\_concat\_through\_map\_dot\_product\_last\_axis(), lic\_array\_vector\_matrix(),



```
(concat
  (elementwise-a
    bsg-systolic
    (slice v-32
    (slice t-32-
    bsg-systolic<sup>.</sup>
    (slice v-32
    (slice t-32-
  (elementwise-a
    bsg-systolic<sup>,</sup>
    (slice v-32
    (slice t-32-
    bsg-systolic<sup>,</sup>
    (slice v-32
    (slice t-32-
  0)"
.parse::<Pattern
.unwrap()
.search_eclass(&runner.egraph, id)
.expect("Did not find expected program");
```

#### (bsg-systolic-array 32 32 v-32 t-32-32)" .parse::<Pattern<\_>>() .unwrap() .search\_eclass(&runner.egraph, id) .expect("Did not find expected program");

```
"(concat
  (bsg-systolic-a
   v-32
   (slice t-32-32
  bsg-systolic-a
   v-32
   (slice t-32-32
  0)"
.parse::<Pattern<
.unwrap()
.search_eclass(&ru
.expect("Did not
```

rray 32 16	
0 32 0 16)	
rray 32 16	
0 32 16 32)	
_>>()	
unner.egraph, find expected	<pre>id) program");</pre>

"(elementwise-add bsg-systolic-array 16 32 (slice v-32 0 16) (slice t-32-32 0 16 0 32) bsg-systolic-array 16 32 (slice v-32 16 32) (slice t-32-32 16 32 0 32) - 11 .parse::<Pattern<\_>>() .unwrap() .search\_eclass(&runner.egraph, id)



```
(concat
  bsg-systolic-array 32 16
  v-32
   (slice t-32-32 0 32 0 16)
  (elementwise-add
   bsg-systolic-array 16 16
    (slice v-32 0 16)
    (slice t-32-32 0 16 16 32)
   bsg-systolic-array 16 16
    (slice v-32 16 32)
 0)"
.parse::<Pattern<_>>()
.unwrap()
.search_eclass(&runner.egraph, id)
```

## .expect("Did not find expected program");

(slice t-32-32 16 32 16 32)

#### • Initial success of this first experiment is promising!

- Initial success of this first experiment is promising!
- Many questions left to be answered

ent is promising! 1

- Initial success of this first experiment is promising!
- Many questions left to be answered
  - Will this representation of tensor programs continue to work?

- Initial success of this first experiment is promising!
- Many questions left to be answered

  - Will this representation of tensor programs continue to work? - How do we pull a "compiler" out of an egraph?

- Initial success of this first experiment is promising!
- Many questions left to be answered
  - Will this representation of tensor programs continue to work?
  - How do we pull a "compiler" out of an egraph?
  - How to implement a complex, multi-objective extraction process?

- Initial success of this first experiment is promising!
- Many questions left to be answered
  - Will this representation of tensor programs continue to work?
  - How do we pull a "compiler" out of an egraph?

  - How to implement a complex, multi-objective extraction process? - How to represent hardware sharing?





# Service Service Sauper Sauper Service Service







PAUL G. SCHOOL 

# Thank you!

https://github.com/gussmith23/glenside https://justg.us