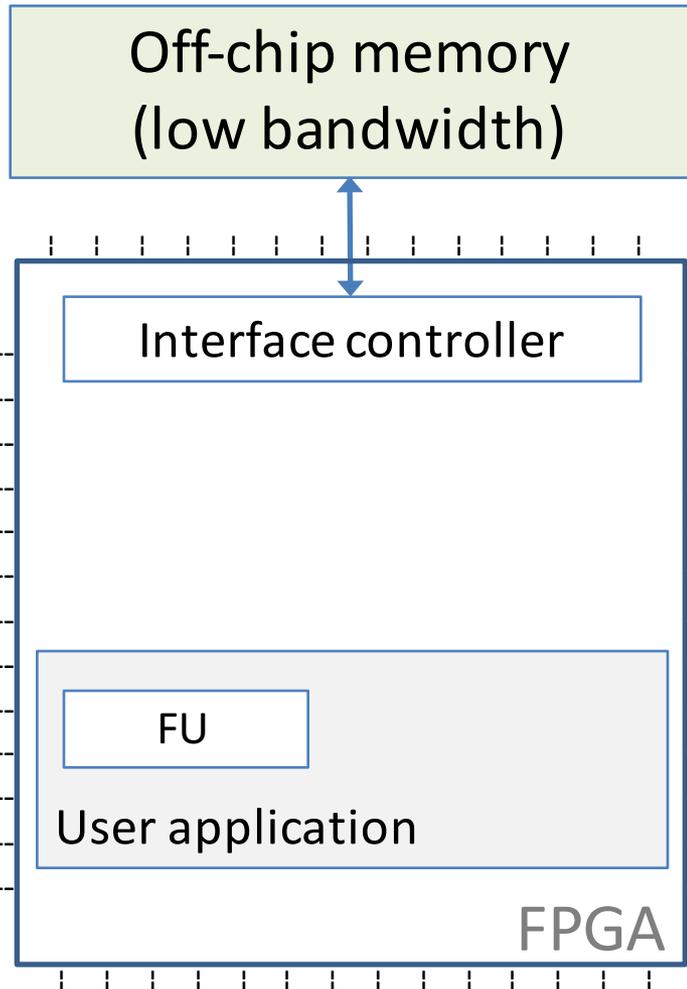


# Custom-Sized Caches in Application-Specific Memory Hierarchies

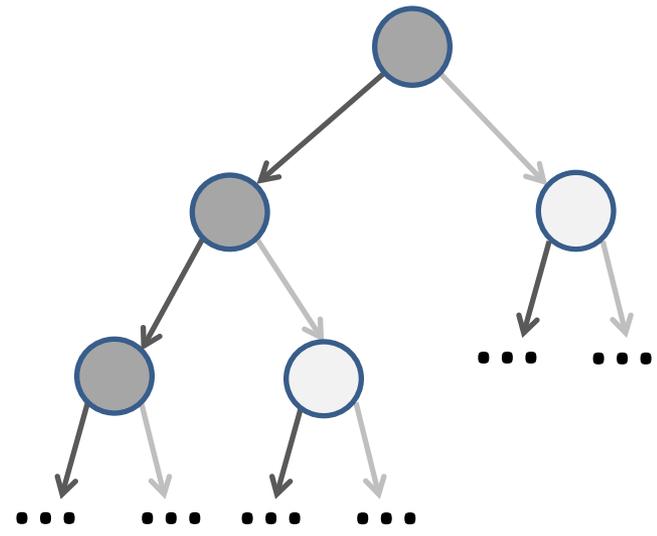
Felix Winterstein, Kermin Fleming, Hsin-Jung Yang,  
John Wickerson and George Constantinides

9 December 2015

# Recap: FPGA'15

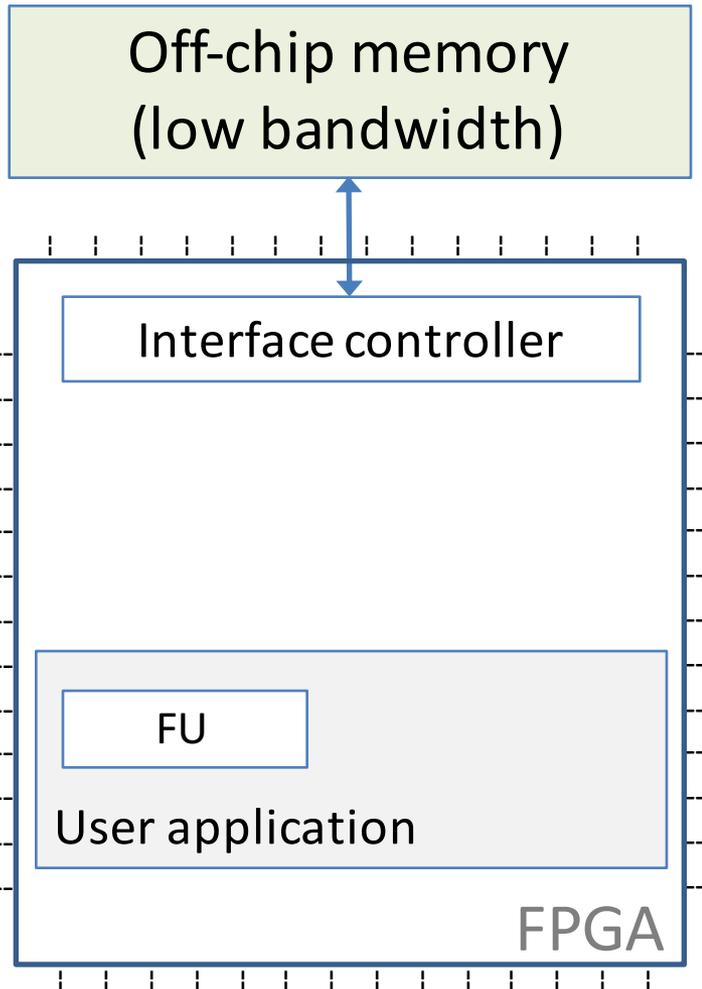


HLS  
←

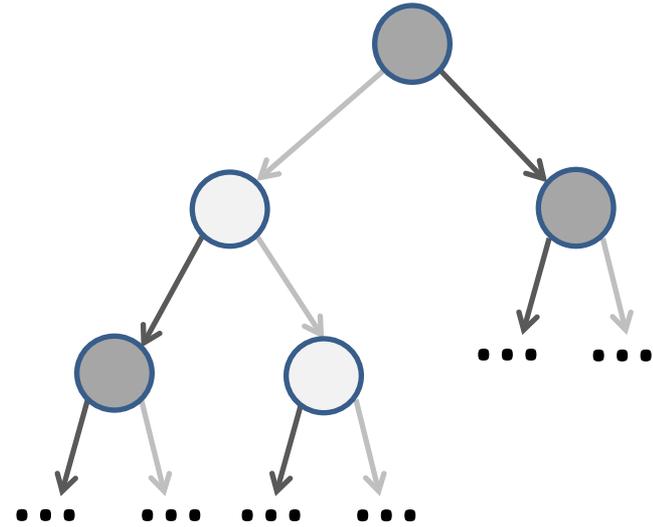


```
void reflectTree (typeB *root)
  typeA *s = new typeA;
  s->u = root;
  s->n = 0;
  while s!=0 do
    typeB *u = s->u;
    ...
```

# Recap: FPGA'15

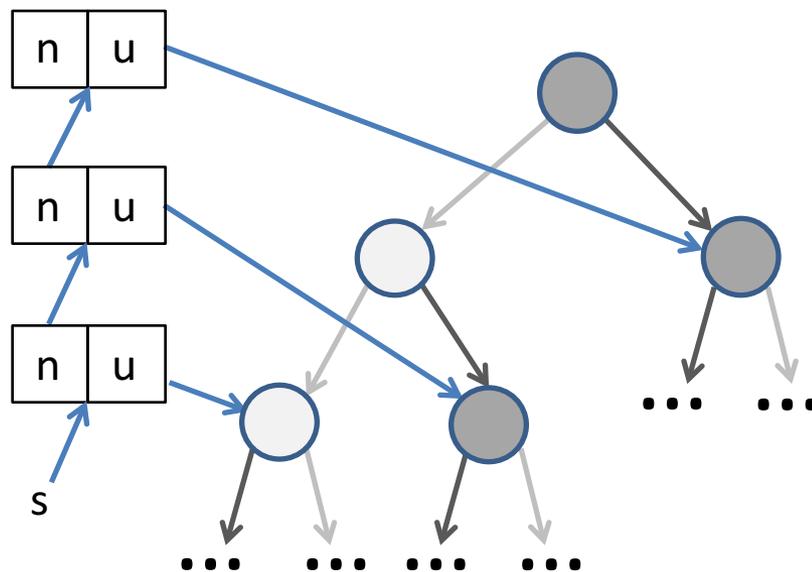
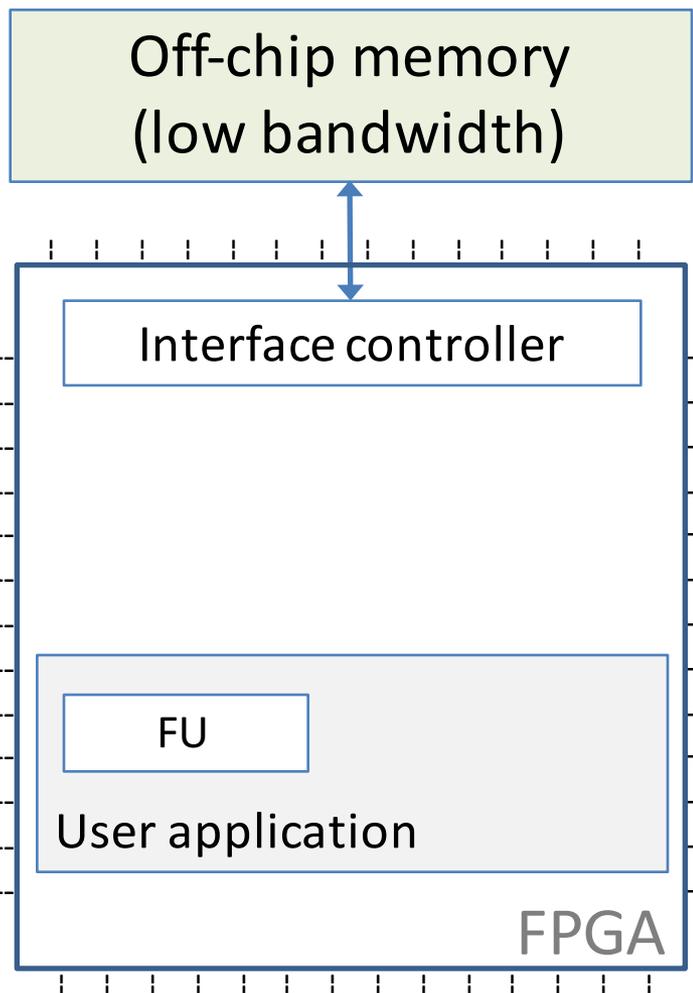


HLS  
←



```
void reflectTree (typeB *root)
  typeA *s = new typeA;
  s->u = root;
  s->n = 0;
  while s!=0 do
    typeB *u = s->u;
    ...
```

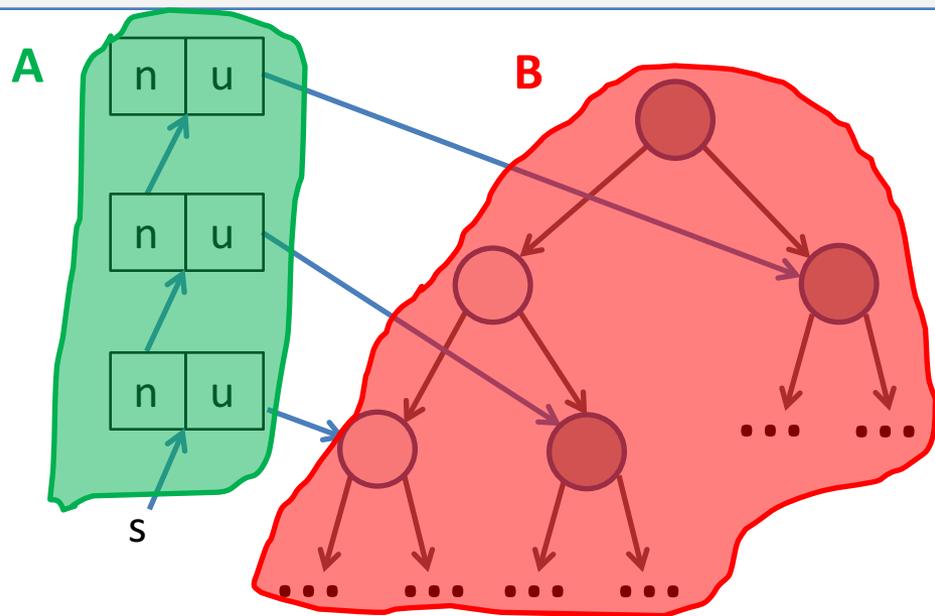
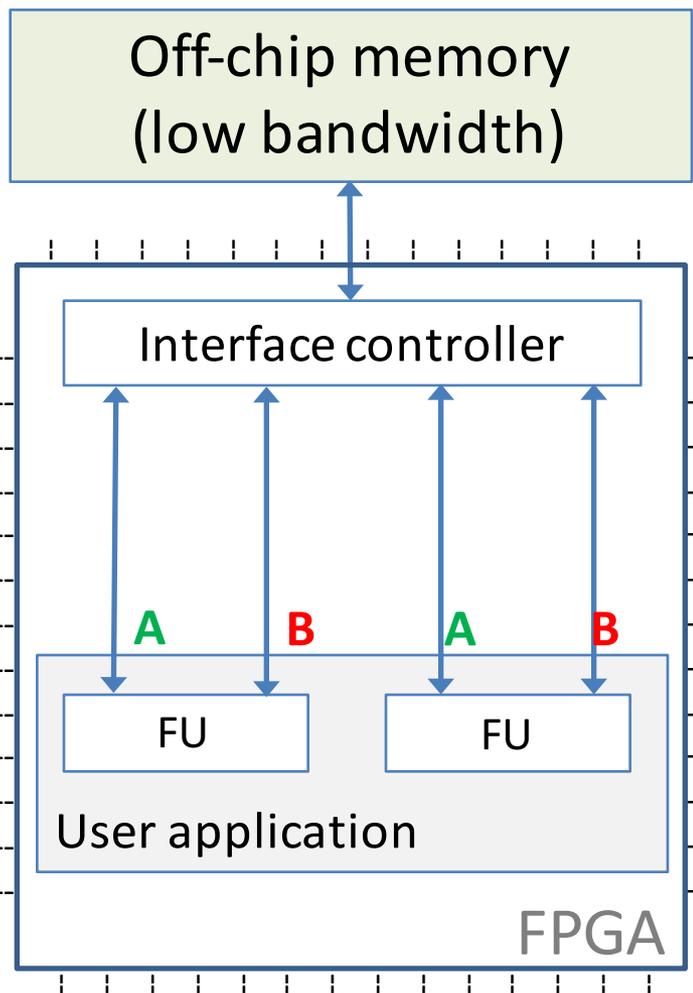
# Recap: FPGA'15



HLS  
←

```
void reflectTree (typeB *root)
typeA *s = new typeA;
s->u = root;
s->n = 0;
while s!=0 do
    typeB *u = s->u;
    ...
```

# Recap: FPGA'15

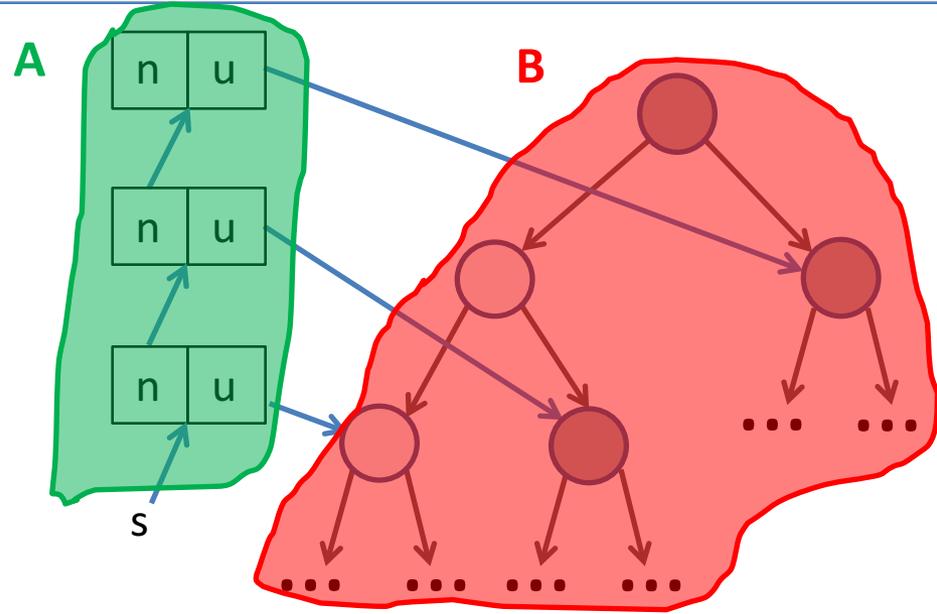
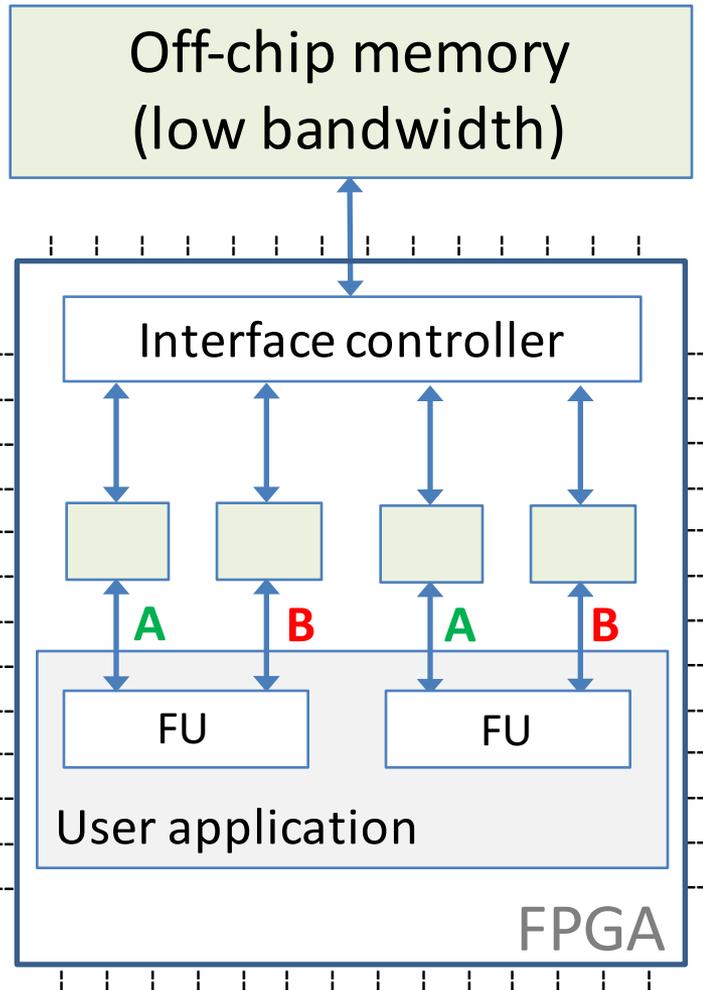


HLS



```
void reflectTree (typeB *root)
typeA *s = new typeA;
s->u = root;
s->n = 0;
while s!=0 do
    typeB *u = s->u;
    ...
```

# Recap: FPGA'15

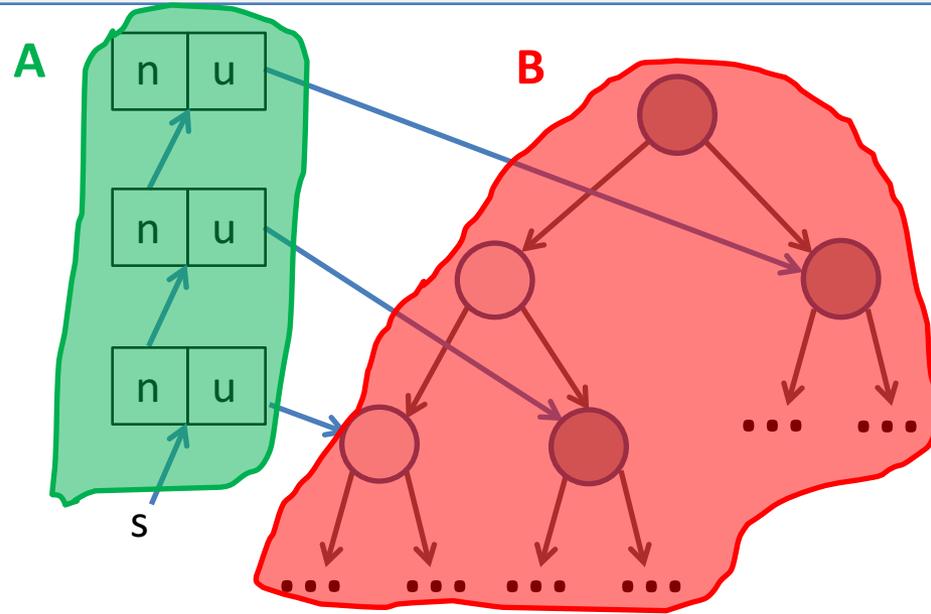
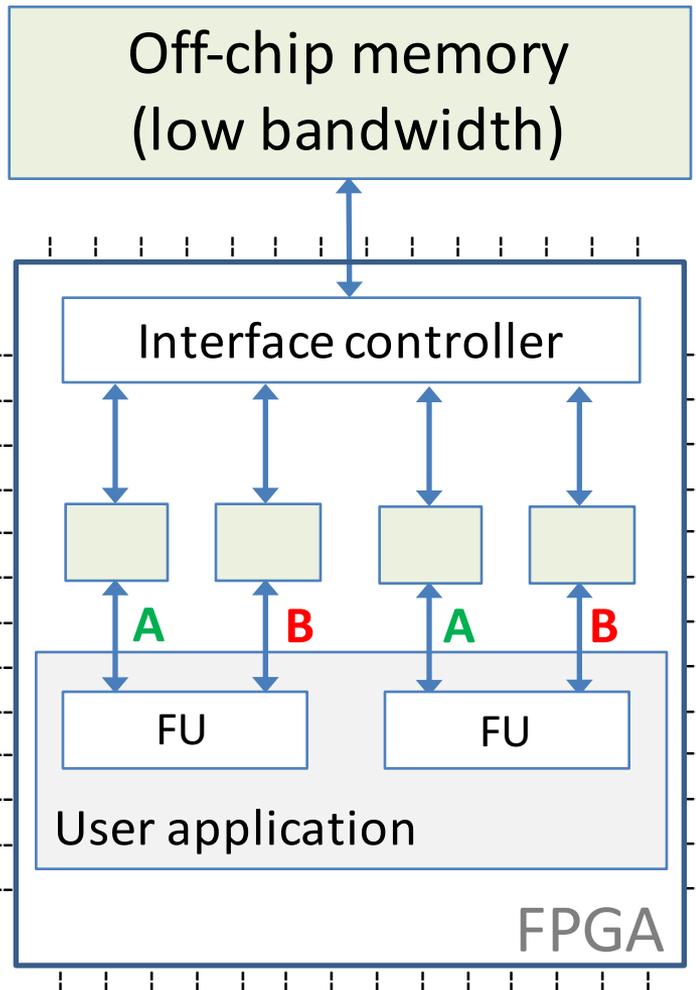


HLS



```
void reflectTree (typeB *root)
typeA *s = new typeA;
s->u = root;
s->n = 0;
while s!=0 do
    typeB *u = s->u;
    ...
```

# Goal



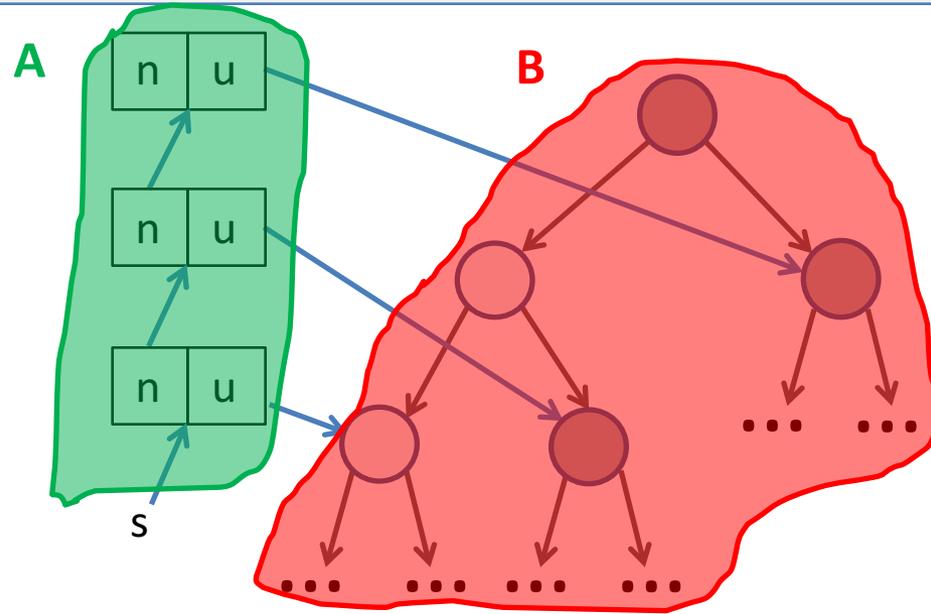
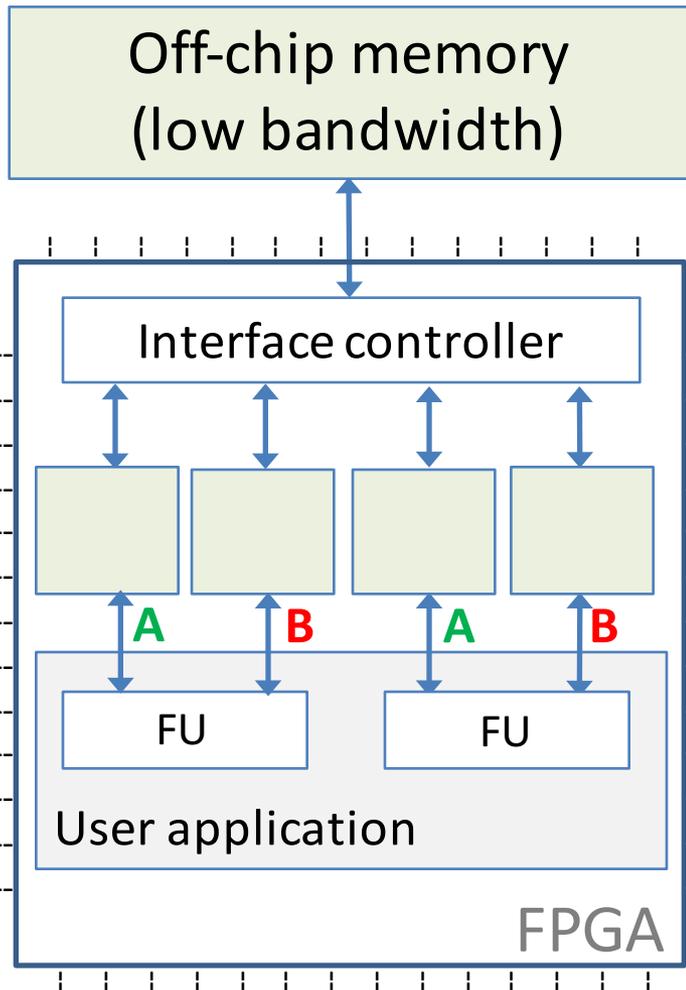
HLS



Custom  
cache  
sizing

```
void reflectTree (typeB *root)
typeA *s = new typeA;
s->u = root;
s->n = 0;
while s!=0 do
    typeB *u = s->u;
    ...
```

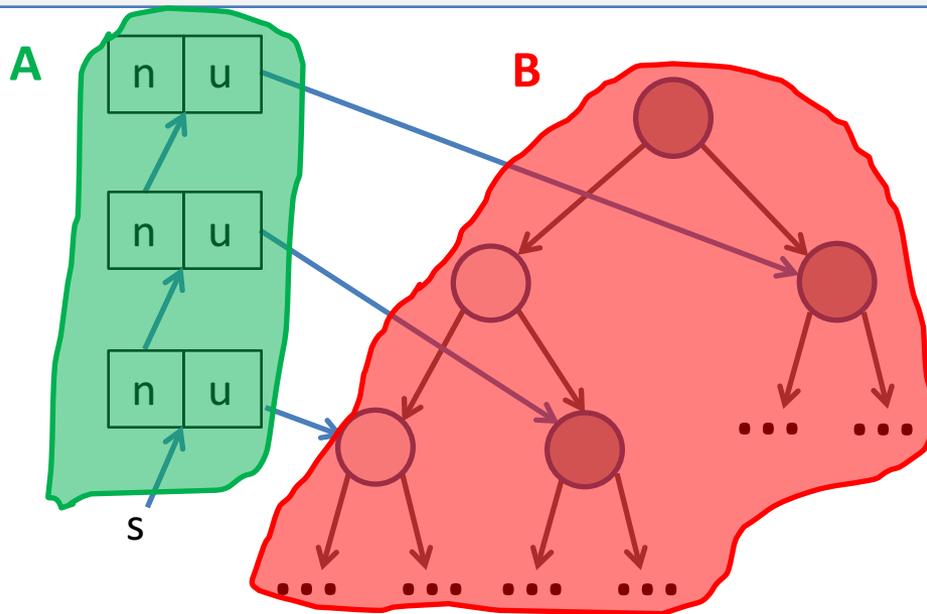
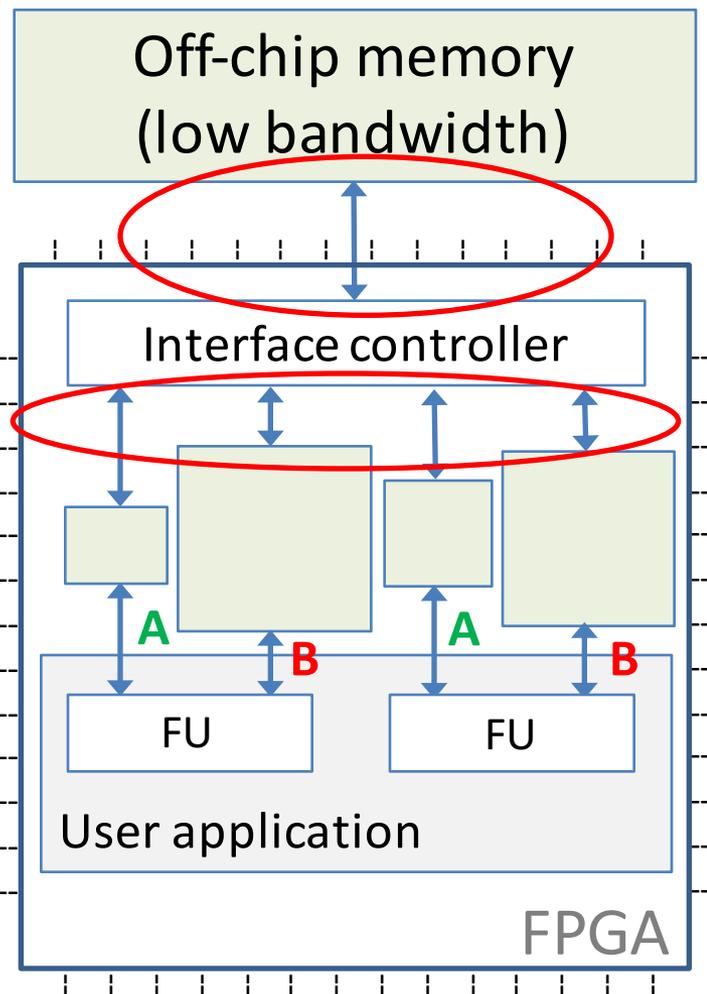
# Goal



**HLS**  
←  
Custom  
cache  
sizing

```
void reflectTree (typeB *root)
typeA *s = new typeA;
s->u = root;
s->n = 0;
while s!=0 do
    typeB *u = s->u;
    ...
```

# Goal



**HLS**  
←  
Custom  
cache  
sizing

```
void reflectTree (typeB *root)
typeA *s = new typeA;
s->u = root;
s->n = 0;
while s!=0 do
    typeB *u = s->u;
    ...
```

# Executive summary

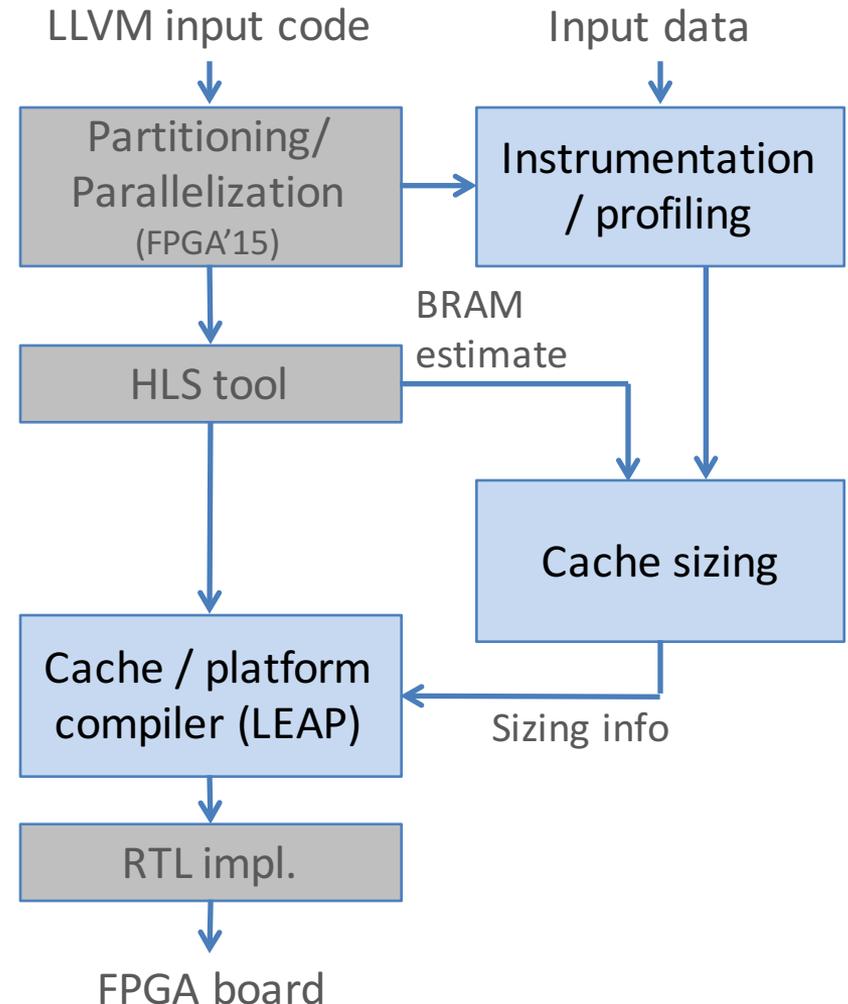
- Maximize aggregate hit rate of multi-cache system
- Scale up caches using “left-over” BRAM
- Assign individual size to each cache

Compared to uncached app:

- **3.6x** latency reduction (average)
- **2.6x** energy reduction (average)

Compared to default size:

- **1.4x** latency savings from variable sizing (average)



# Remainder of this talk

1. Estimate left-over BRAM
2. Variably-sized caches
3. Implementation
4. Evaluation

# Estimate left-over BRAM

- BRAM estimate from HLS tool + known offsets
- ReflectTree example on a Virtex 7 FPGA:

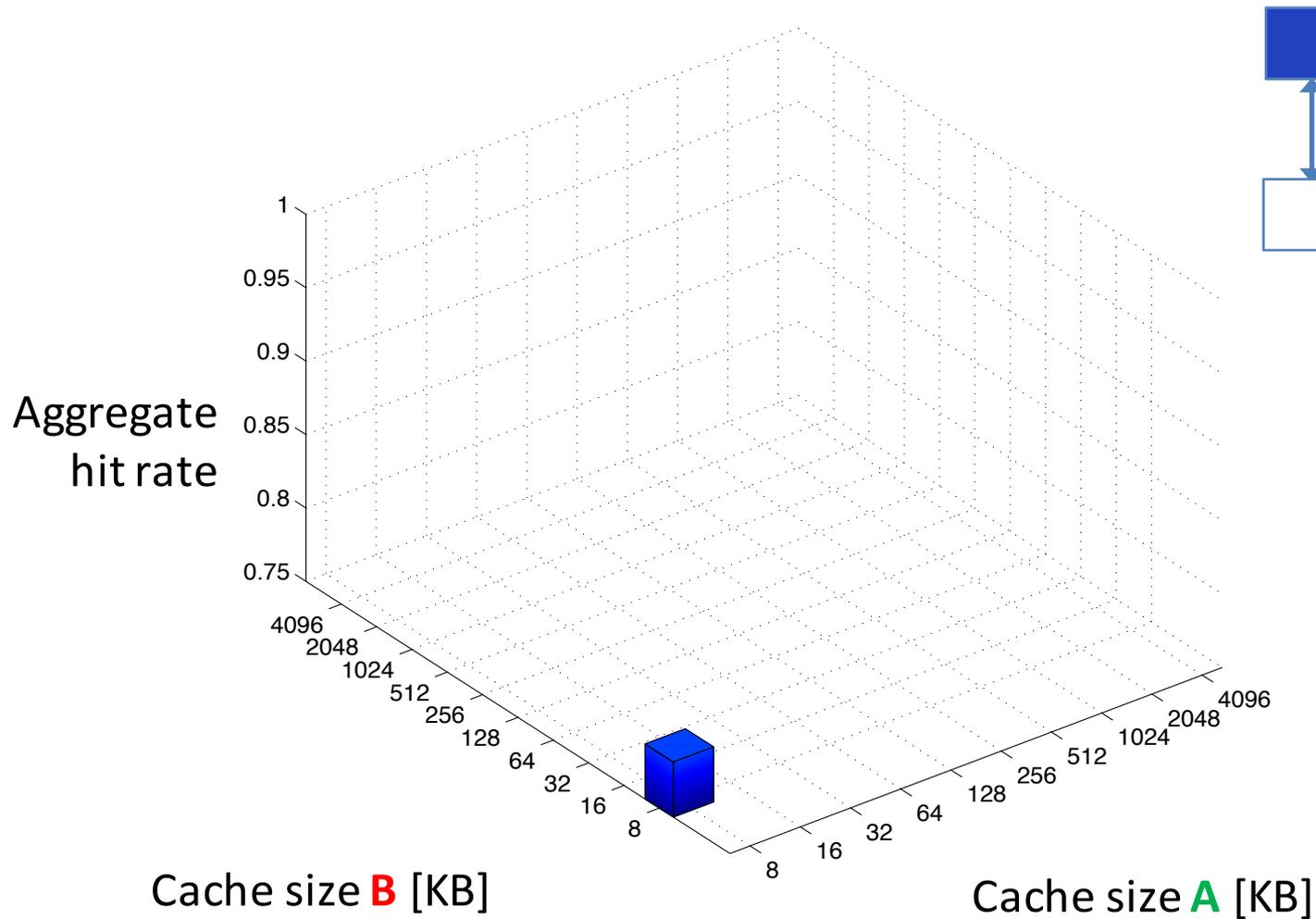
Component	Estimate [RAM blocks]	Post-PAR [RAM blocks]
HLS core	208	158
Memory interfaces	40	40
LEAP platform (without scratchpads, fixed offset)	50.5	50.5
<b>Unused left-over BRAM (Virtex 7)</b>	<b>628.5</b>	<b>678.5</b>

~2200 KB can be repurposed for cache construction

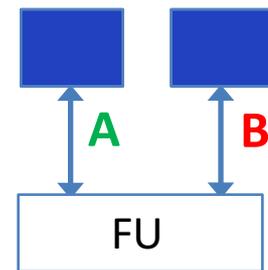
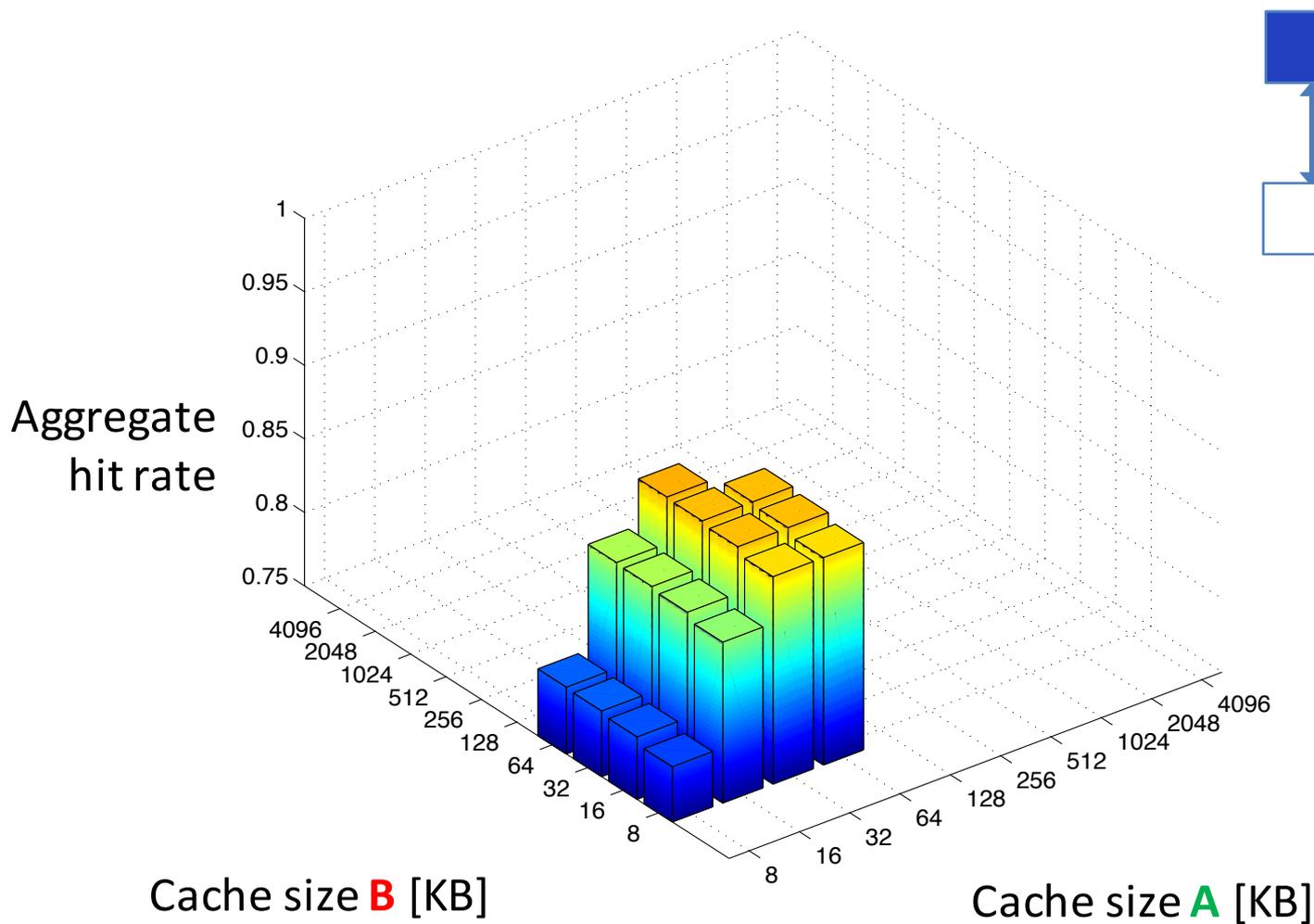
# Outline

1. Estimate left-over BRAM
2. Variably-sized caches
3. Implementation
4. Evaluation

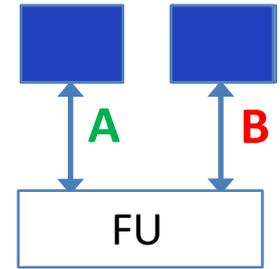
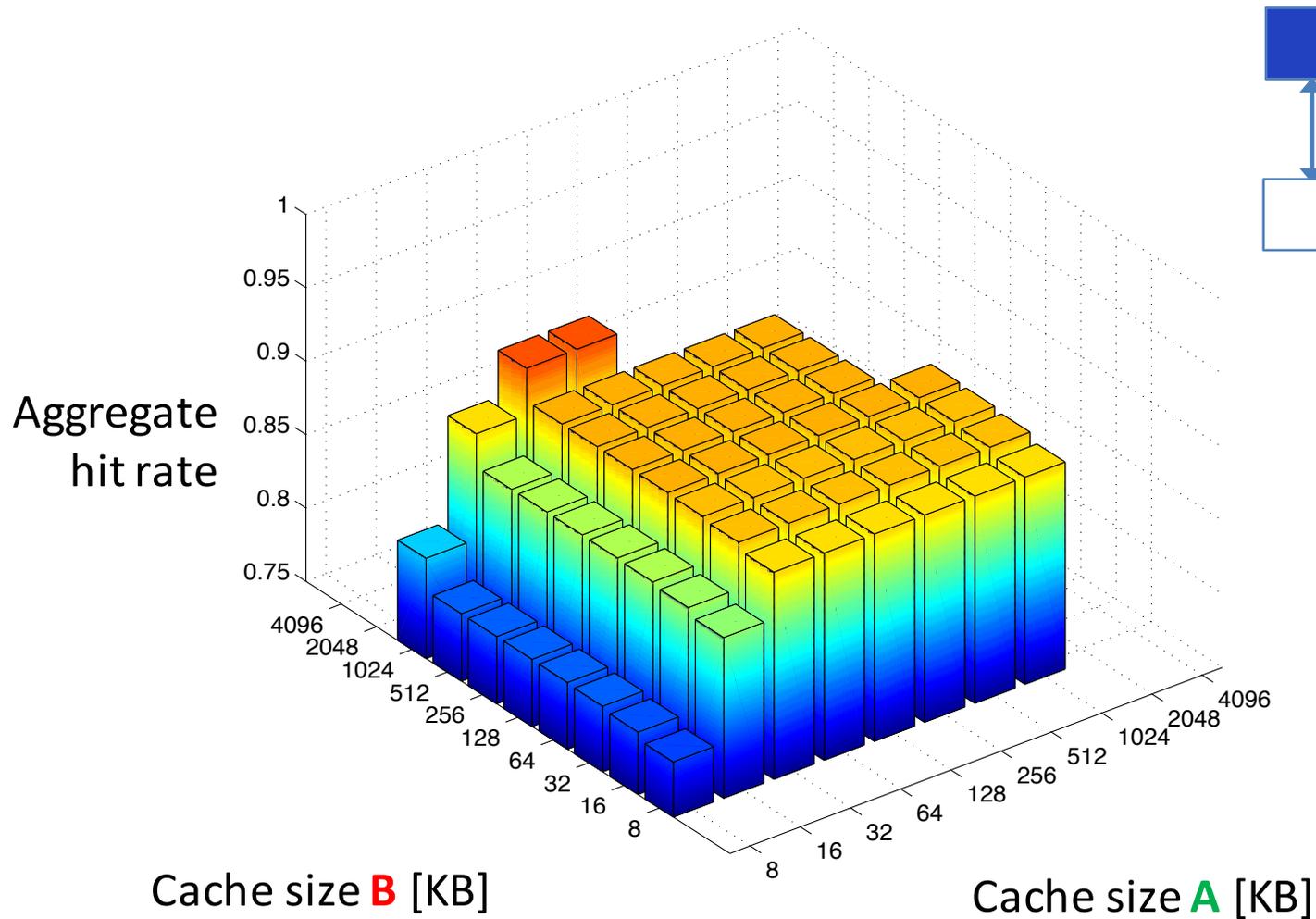
# Variably-sized caches



# Variably-sized caches

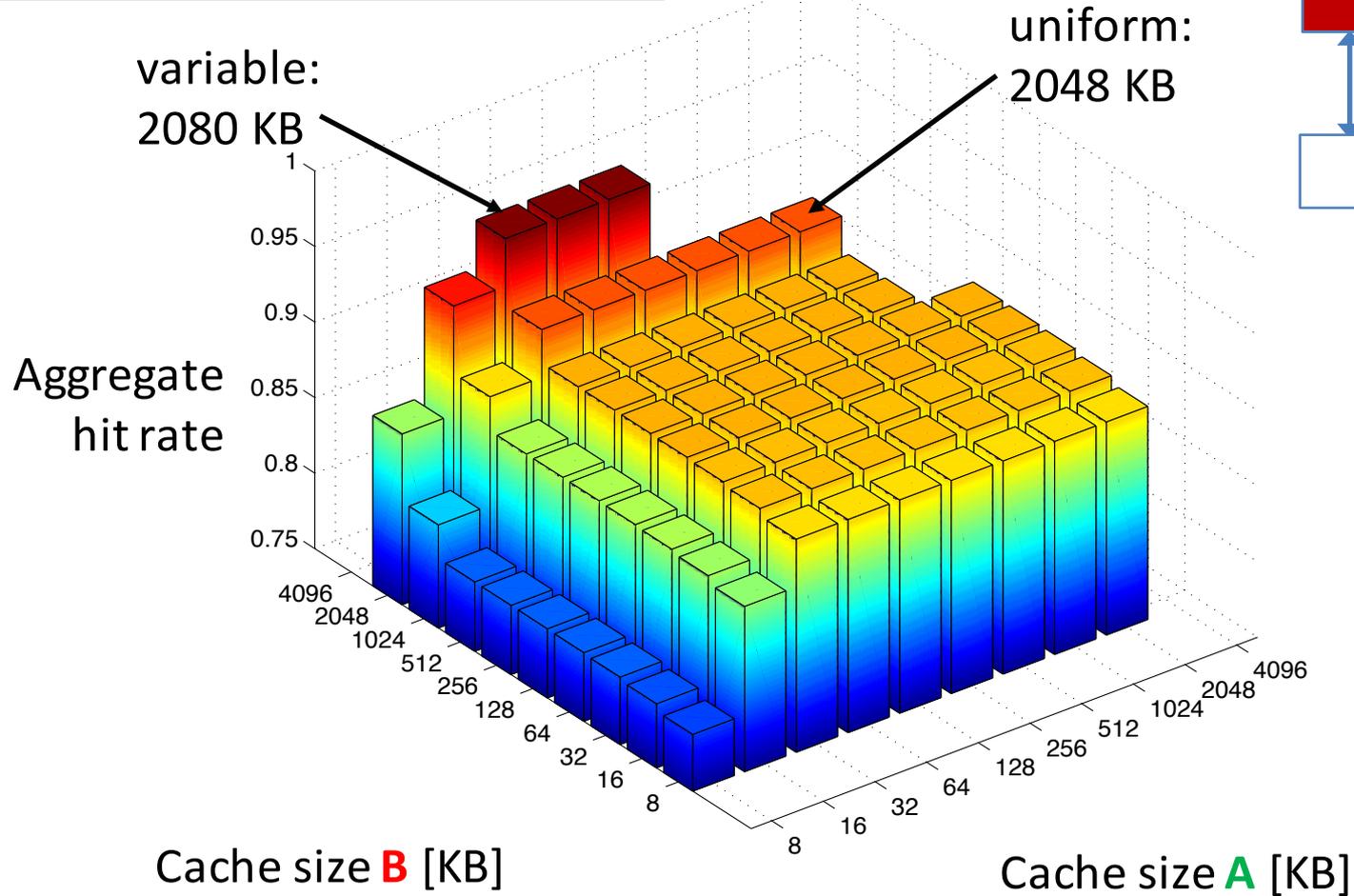


# Variably-sized caches

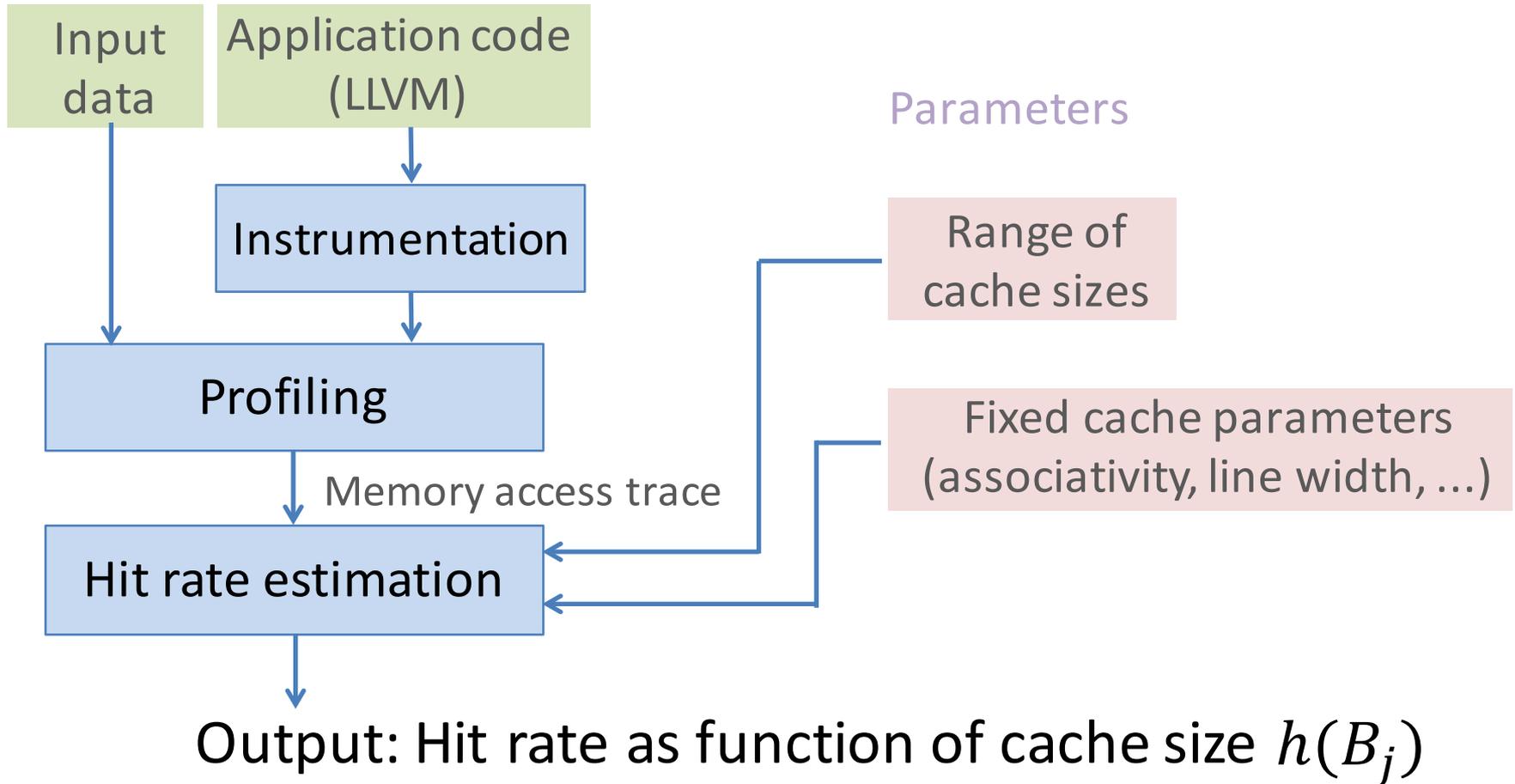


# Variably-sized caches

Resource constraint: 2200 KB BRAM



# Cache performance model



# Optimization strategy

Task:

- Find “optimal” size assignment for  $K$  caches
- Choose among  $N$  sizes for each cache

Cast into Multiple-Choice Knapsack Problem:

maximize  $\sum_{i=1}^K \sum_{j=1}^N h_i(B_j) x_{ij}$  hit rate for cache size  $B_j$

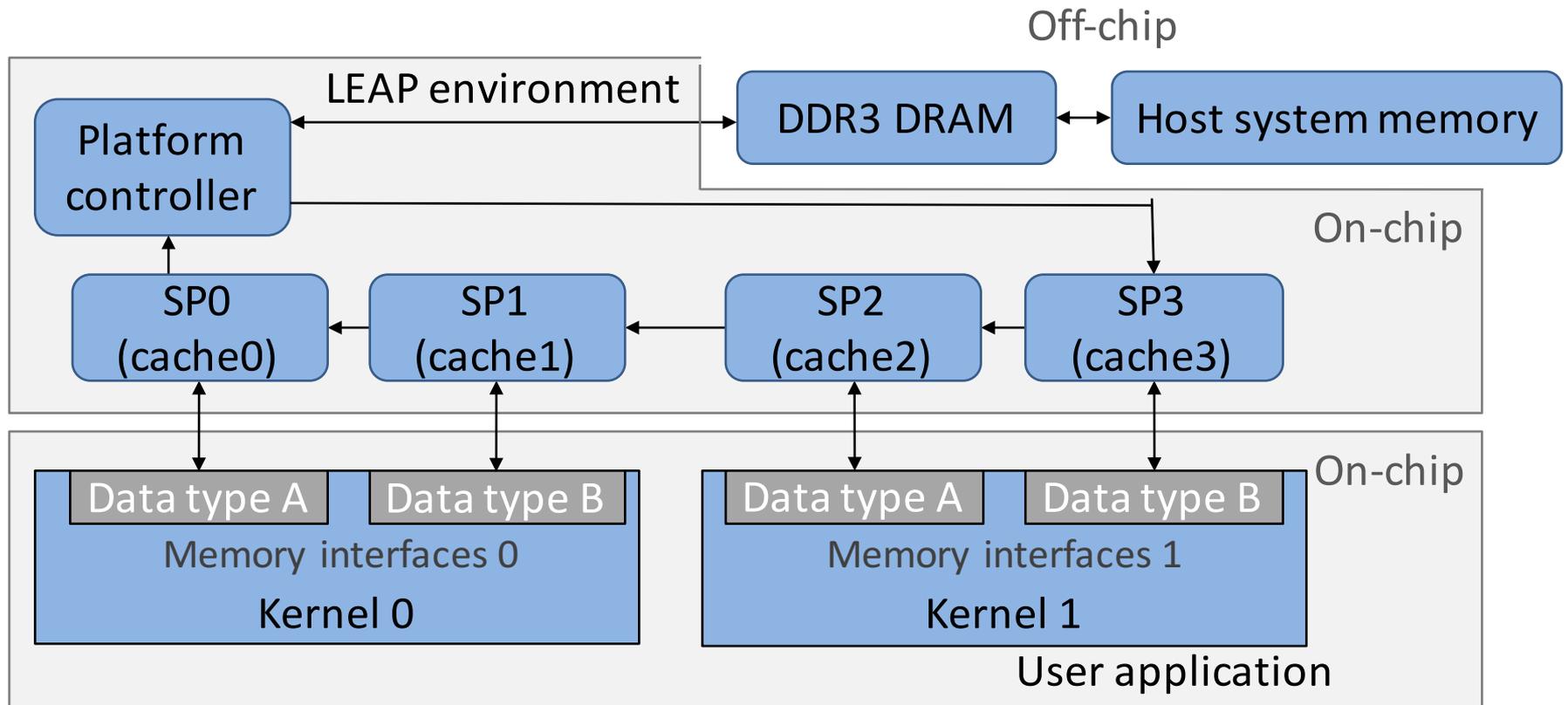
subject to  $\sum_{i=1}^K \sum_{j=1}^N BRAM_i(B_j) x_{ij} \leq C$  Global BRAM constraint

and  $\sum_{i=1}^K x_{ij} = 1, \quad i = 1 \dots K$  Select each cache exactly once

# Outline

1. Estimate left-over BRAM
2. Variably-sized caches
- 3. Implementation**
4. Evaluation

# Implementation



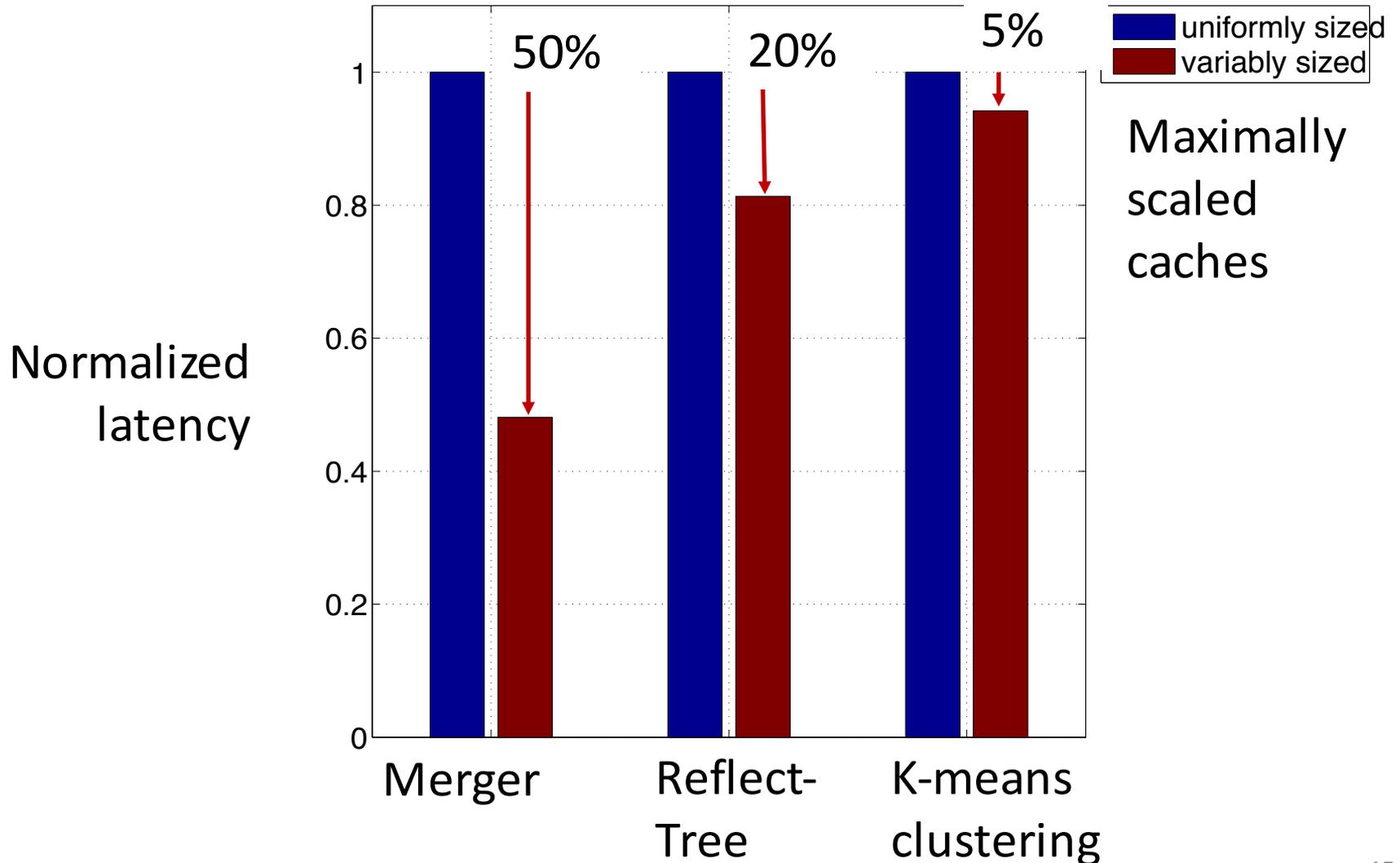
# Outline

1. Estimate left-over BRAM
2. Variably-sized caches
3. Implementation
4. Evaluation

# Results

	LUT	FF	BRAM	Latency	
<b>1 Merger (8 scratchpads)</b>					
No caches	67k	67k	586.5	7.7 ms	 x5.0
Variably-sized, scaled caches	92k	88k	874.5	1.5 ms	
<b>2 ReflectTree (4 scratchpads)</b>					
No caches	74k	77k	248.5	345.6 ms	 x3.4
Variably-sized, scaled caches	85k	83k	944.5	100.3 ms	
<b>3 K-means clustering (6 scratchpads)</b>					
No caches	91k	91k	347.5	598.7 ms	 x2.7
Variably-sized, scaled caches	107k	103k	829.5	221.8 ms	

# Benefit of custom sizing



# Energy

		Energy (FPGA+DRAM)	
<b>1</b>	<b>Merger (8 scratchpads)</b>		
	No caches	22.0 mJ	 x3.9
	Variably-sized, scaled caches	5.6 mJ	
<b>2</b>	<b>ReflectTree (4 scratchpads)</b>		
	No caches	1060.5 mJ	 x2.2
	Variably-sized, scaled caches	477.6 mJ	
<b>3</b>	<b>K-means clustering (6 scratchpads)</b>		
	No caches	1976.9 mJ	 x2.0
	Variably-sized, scaled caches	973.5 mJ	

# Conclusion

- Automatic cache size scaling in a multi-cache system
  - Using spare BRAM resources
  - Profiling-based hit rate estimation
  - Find the resource allocation that maximizes aggregate hit rate
- Future work
  - Model of a coherency network
  - Optimal partitioning of the memory controller network
  - High-level energy estimation and explicit minimization of energy consumption

# Questions?

---

**Thank you.**