

# Extending the Roofline Model: Bottleneck Analysis with Microarchitectural Constraints

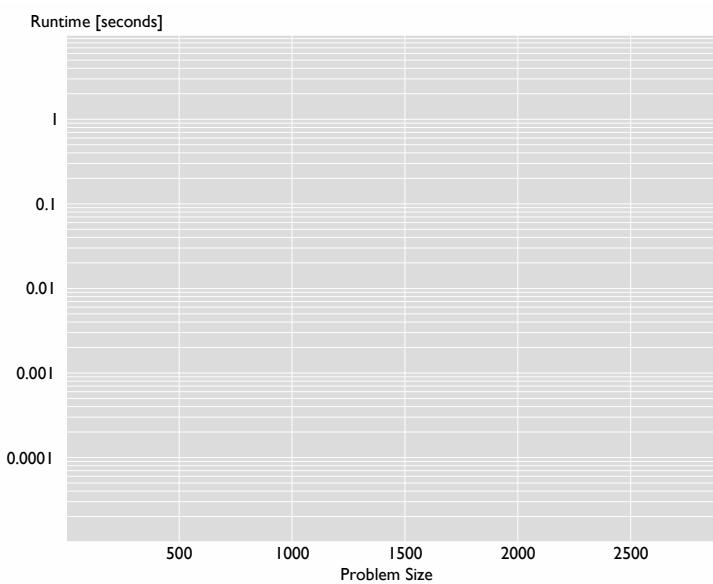
Victoria Caparrós Cabezas  
Markus Püschel



*Part of PhD work*

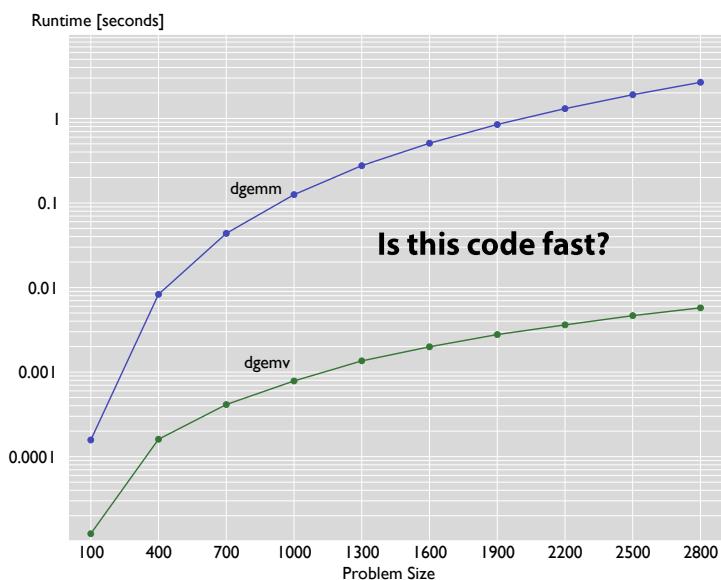
Computer Science  
**ETH** zürich

## Measuring Runtime



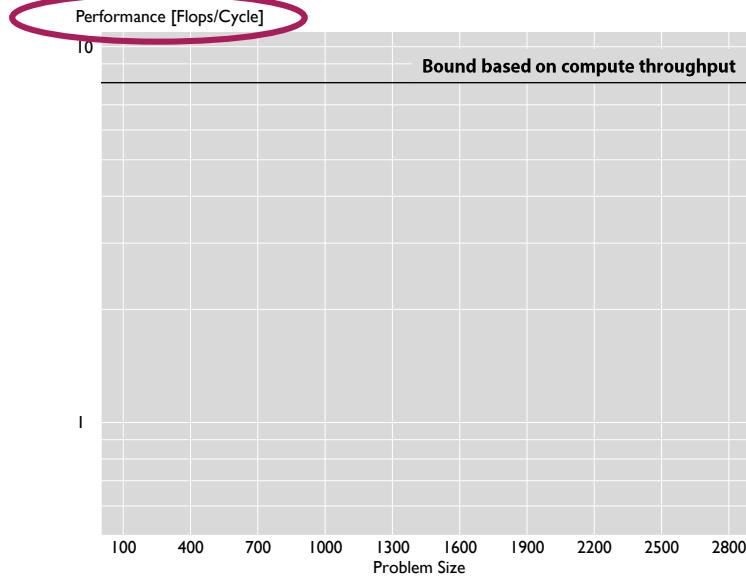
2

## Measuring Runtime



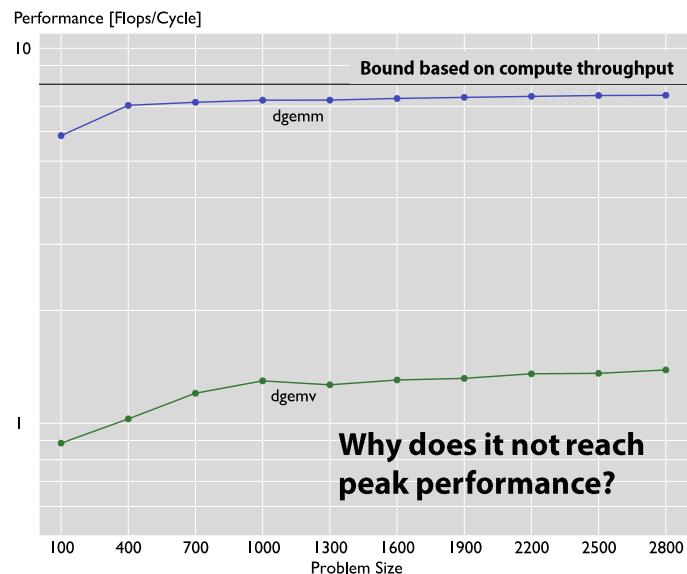
3

## Measuring Performance



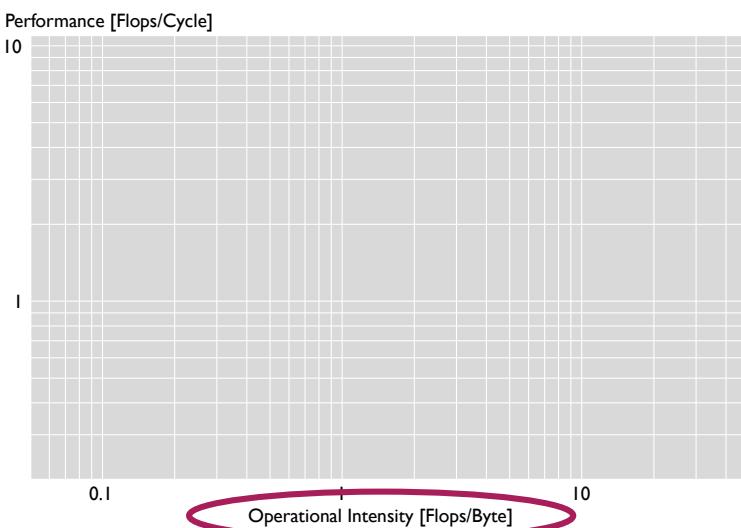
4

## Measuring Performance



5

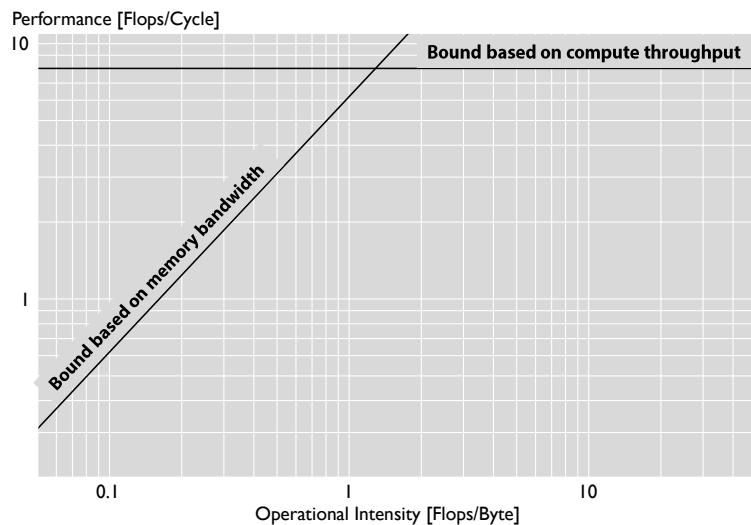
## Roofline Plot



6

Williams et al., "Roofline: An Insightful Visual Performance Model for Multicore",  
Communications of the ACM, 2009

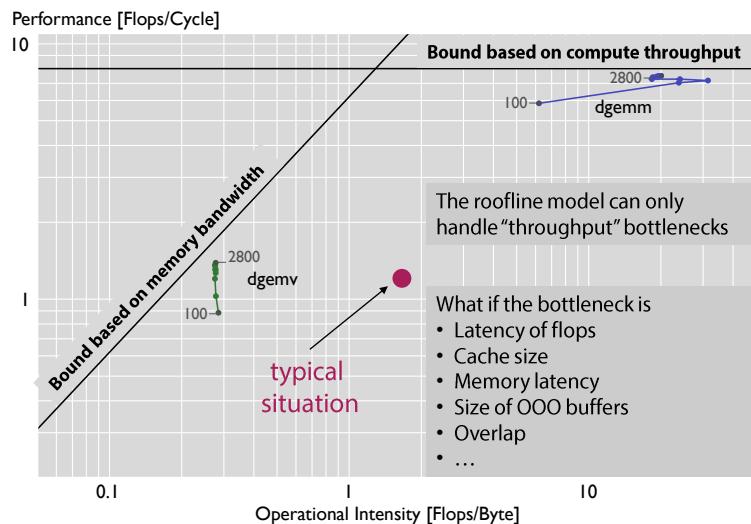
## Roofline Plot



Williams et al. "Roofline: An Insightful Visual Performance Model for Multicore,"  
Communications of the ACM, 2009

7

## Roofline Plot

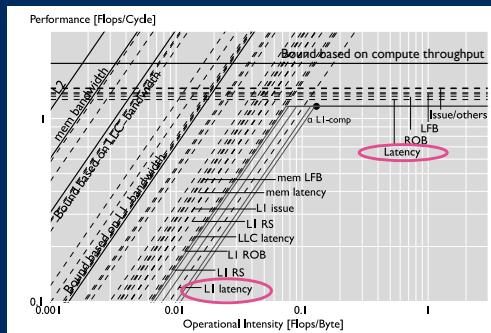


Measurements: Ofenbeck et al. "[Applying the Roofline Model](#)" Proc. IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), 2014, pp. 76-85

8

# Goal

Extended Roofline Model (ERM) to include additional hardware-related bottlenecks



Performance bottlenecks:

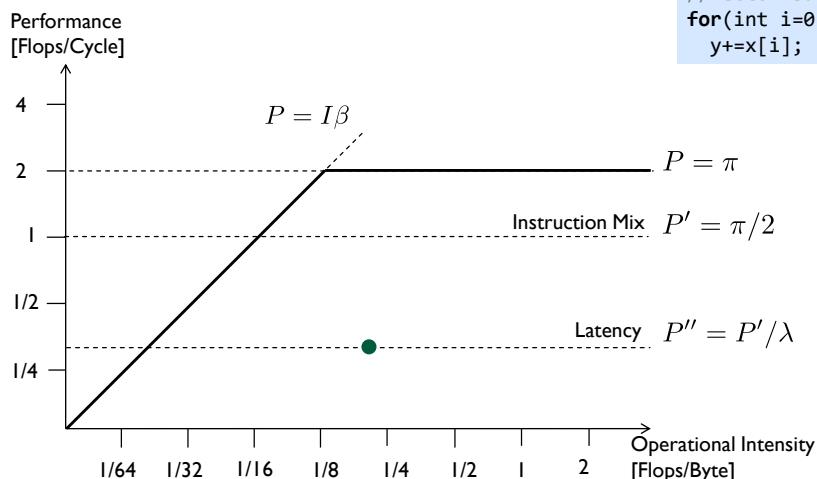
1. Latency of flops
2. Latency of L1 accesses

Current limitation:  
single core (with SIMD) only

## Applications

- Code optimization
- Assess the impact of hardware upgrades

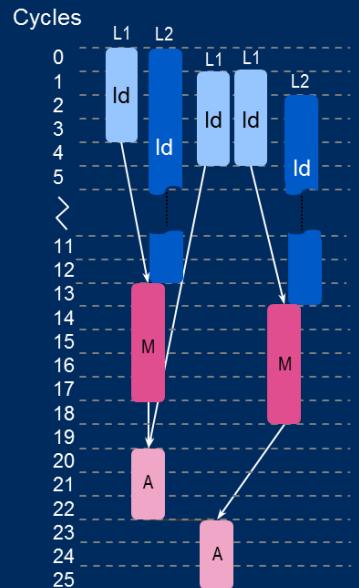
## Example Challenge: Latency



Challenge: What if the latency only affects some of the operations?

## Our solution:

Detailed breakdown of the computation time through a DAG-based performance model



### Input: C Code plus input

```
int main(){
    int n = 1000;
    A=(double*)mm_malloc(...);
    ...
    for(int i=0; i<n; i++)
        for(int j=0; j<n; j++)
            for(int k=0; k<n; k++)
                C[i][j]+=A[i][k]*B[k][j];
}
```

1 clang -emit-llvm -c -O3 ...

LLVM IR  
(no registers, SSA, some vector instructions)

```
%5 = trunc i64 %indvars.lv87 to i32
%idxprom13.us.unr = sext i32 %add12 to i64
%arrayidx14.us.unr = getelementptr %B,
%6 = load double* %arrayidx14.us.unr, align 8,
%mul15.us.unr = fmul double %5, %6
%mul15.us.unr = fmul double %5, %6
...
```

2 LLVM dynamic trace generator  
from Railing et al. [TACO 2015]

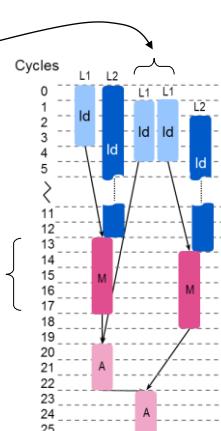
Dynamic execution trace (DAG)

```
38 = load double* %arrayidx18.us, align 8, dbg !385, ibrba !384 (0x3dc7aa8)
36 = load double* %arrayidx17.us, align 8, dbg !385, ibrba !384 (0x3dc7a18)
37 = load double* %arrayidx18.us, align 8, dbg !385, ibrba !384 (0x3dc7aa8)
Result.us = fmul double %36, %37, dbg !385 (0x3dc7d18)
...
```

### Input: Microarchitectural Parameters

Microarchitectural Parameters	Unit	Value
Throughput	flops/cycle	1, 1
$\beta_{mem}, \beta_{L_i}$	doubles/cycle	1, 2/1, 4, 4
$\phi$	instr./cycle	4
Latency	cycles	3, 5
$\lambda_A, \lambda_M$	cycles	100, 4, 12, 30
$\mu_{mem}, \mu_{L_i}$	cycles	
Capacity	bytes	32K, 256K, 20M
$\gamma_{L_i}$	bytes	64
$\varphi$	bytes	
$\gamma_{buf}$	# slots	168, 54, 36, 64, 10

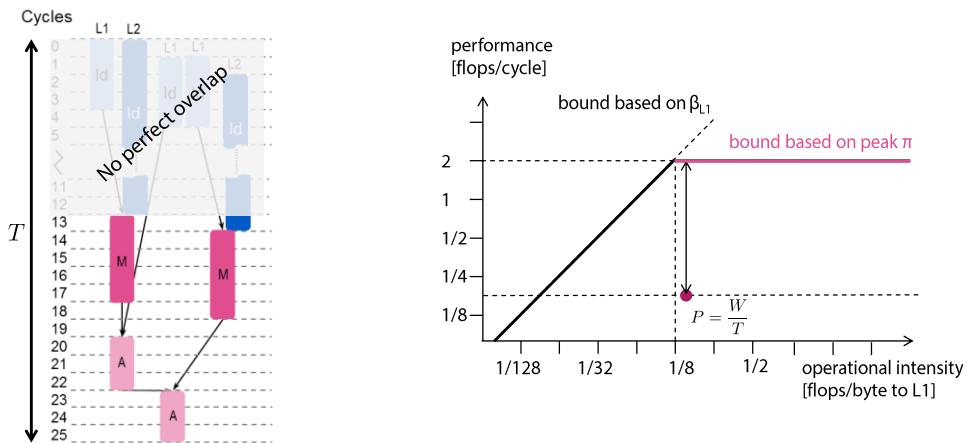
### Output: Scheduled DAG



+ data from its analysis

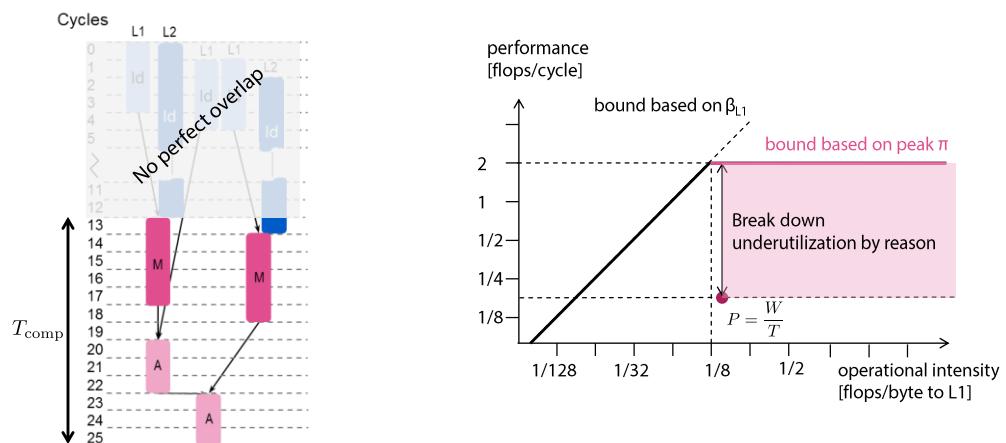
3 Tomasulo's scheduling algorithm  
ERM (Our Tool)

## Utilization-Based Bottleneck Modeling

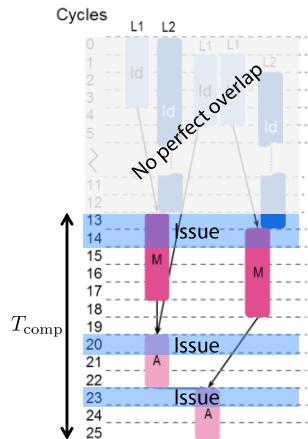


Where are execution cycles spent so that peak performance is not reached?

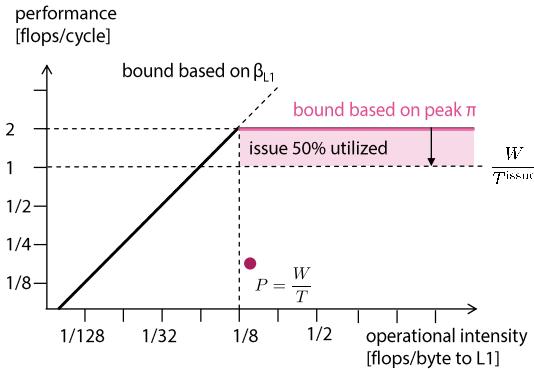
## Utilization-Based Bottleneck Modeling



# Utilization-Based Bottleneck Modeling



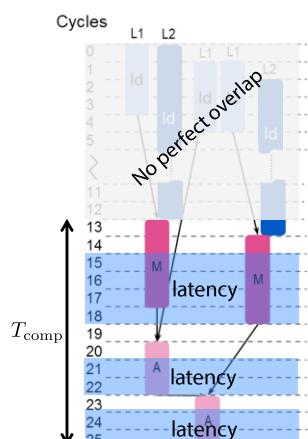
$T^{\text{issue}} = \text{Number of cycles in which nodes are issued}$



$T^{\text{issue}} = 4$

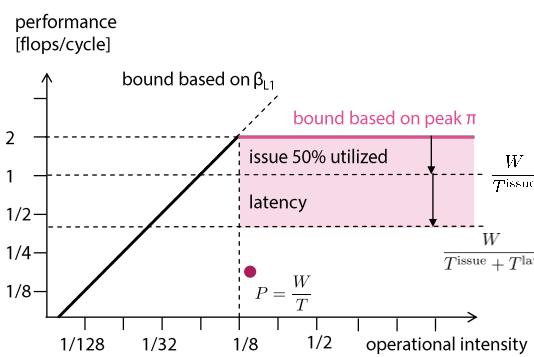
15

# Utilization-Based Bottleneck Modeling



$T^{\text{issue}} = \text{Number of cycles in which nodes are issued}$

$T^{\text{lat}} = \text{Number of cycles in which nodes are executed, but not issued}$



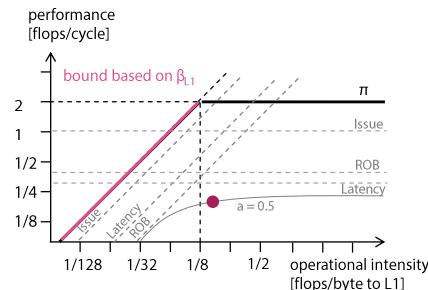
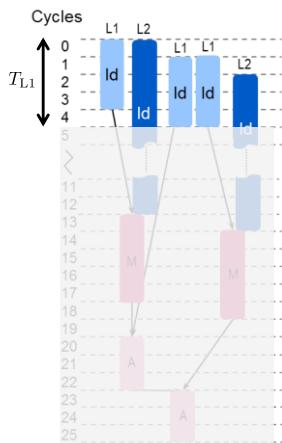
$T^{\text{issue}} = 4$

$T^{\text{lat}} = 8$

16

# Utilization-based Bottleneck Modeling

Similar procedure for each level of the memory hierarchy



Finally, compute overlap (with compute lines)

17

# Merging Roofline Plots for Memory Hierarchy

Original roofline model:

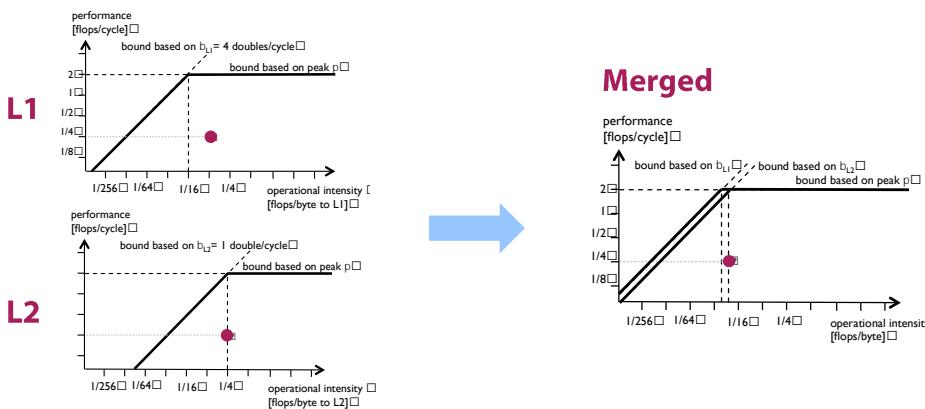
$$I_x = \frac{W}{Q_x}$$

$$P_x \leq \min(I_x \beta_x, \pi)$$

Redefinition of operational intensity:

$$I = \frac{W}{Q} = \frac{W}{Q_{L1} + Q_{L2} + Q_{L3} + Q_{mem}} \quad I_x = I \frac{Q}{Q_x}$$

$$P_x \leq \min(I_x \beta_x, \pi) = \min(I \frac{Q}{Q_x} \beta_x, \pi)$$

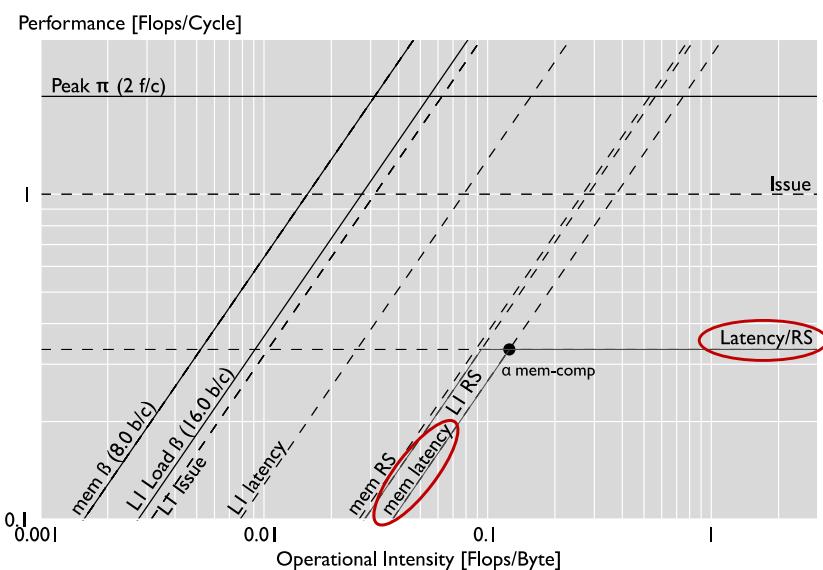


18

## Some Results

19

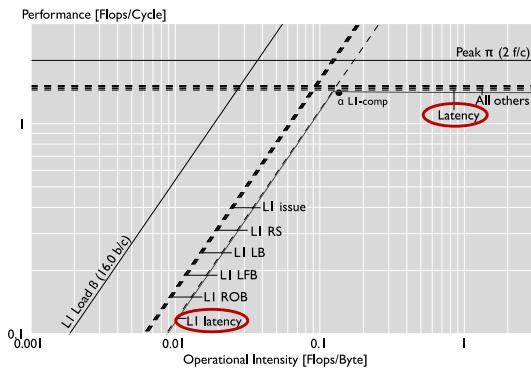
### Vector Sum Reduction (cold cache, $n = 5 \times 10^6$ )



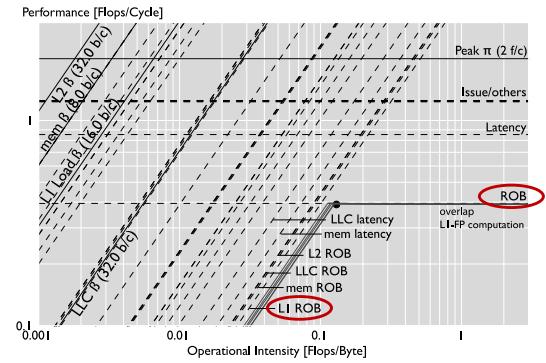
20

## FFT NR (warm cache)— Increasing Size

Size  $2^{10}$



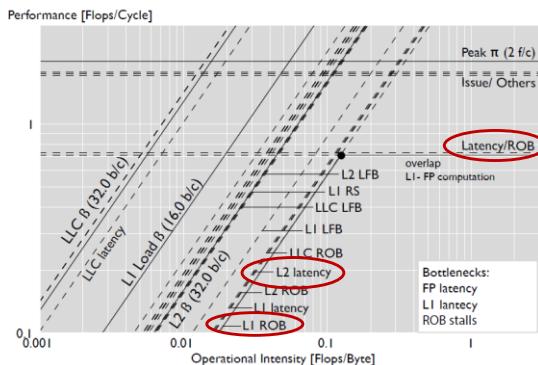
Size  $2^{20}$



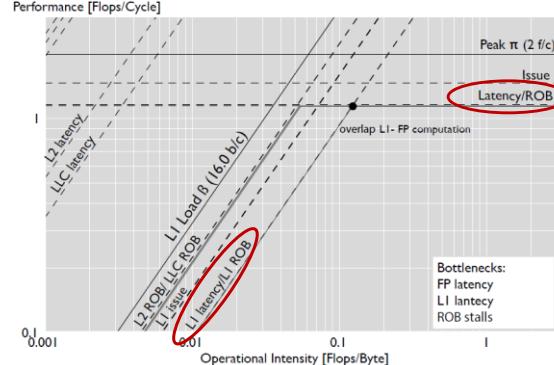
21

## MMM Warm Cache: Triple Loop vs. Blocked ( $n = 500$ )

Triple loop



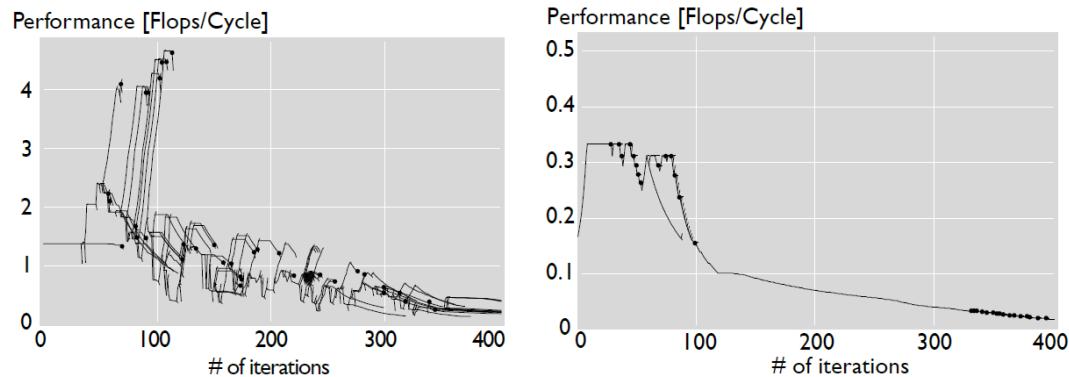
Blocked for L1 ( $b = 50$ )



22

## Followup Work — Balancing Processor for Kernel

Hill climbing algorithm to find balanced processors (= good match SW/HW)

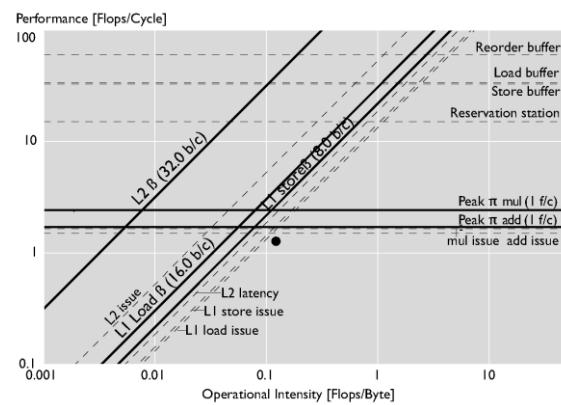


*LL 1: Hydrodynamics computation fragment*

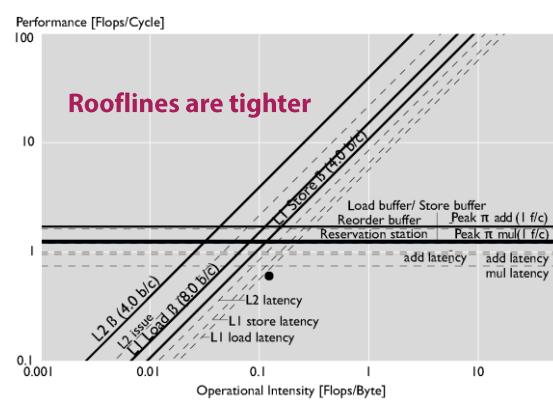
23

## Effect on Roofline Plot (Example)

### On Sandybridge

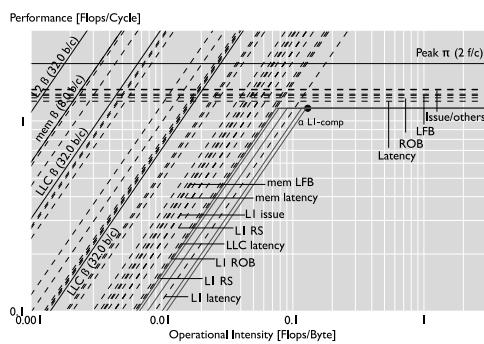


### On a balanced processor



24

# Conclusions



Extended roofline model for computational kernels

- Memory hierarchy
- Throughput and latency
- Size of caches and OOO buffers
- **To date: Single core with SIMD**

**Code:** DAG from LLVM interpreter

**Processor:** 30+ parameters

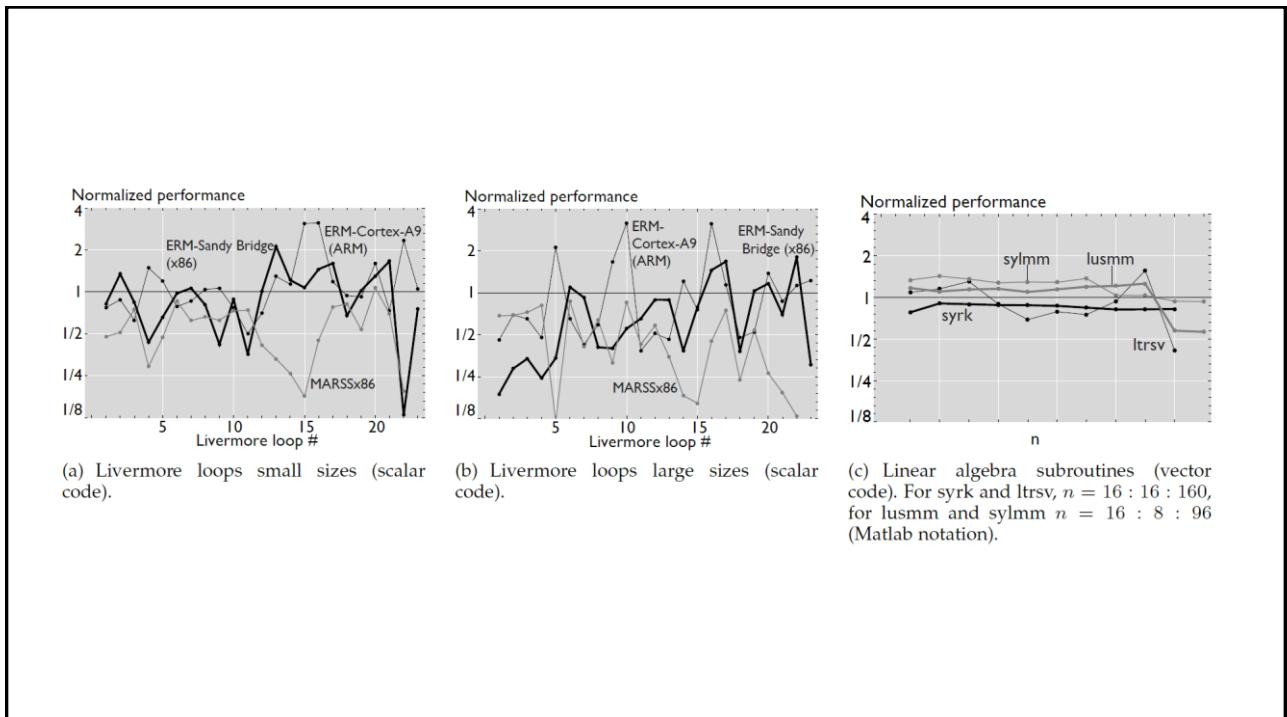
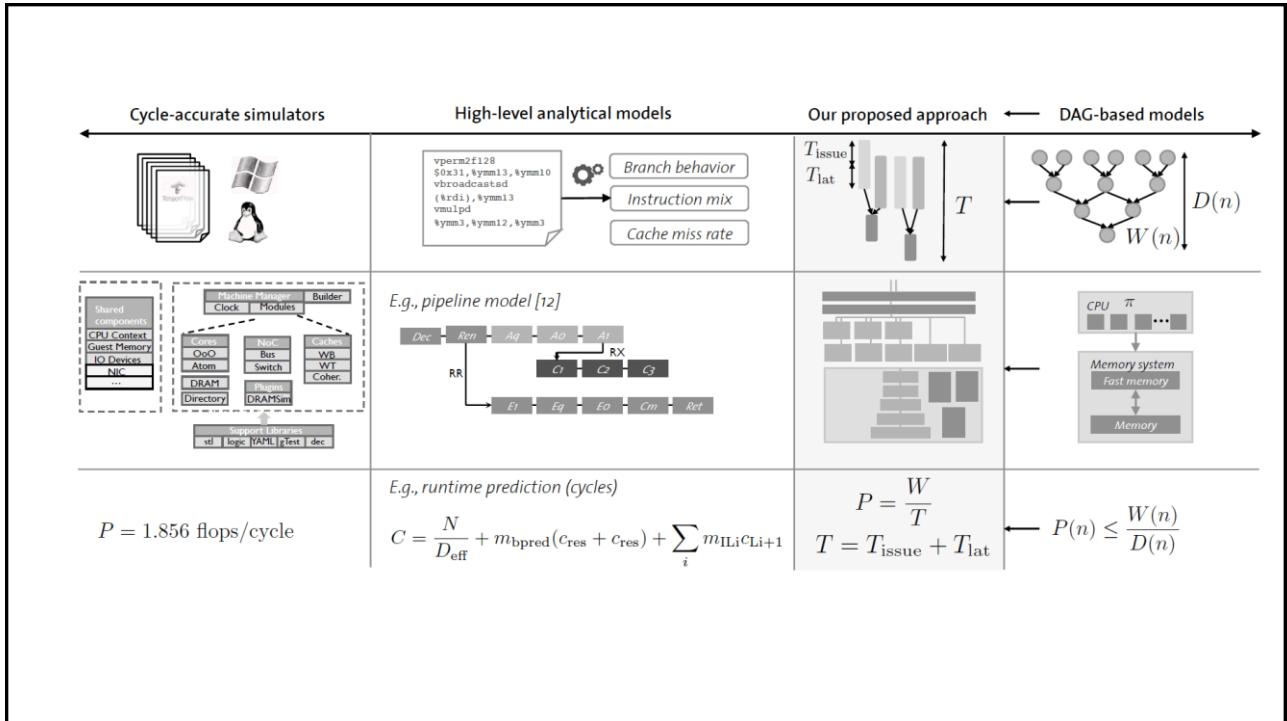
**Basic Idea:** Schedule and analyze DAG

**Tool ERM:**

<https://github.com/caparrov/llvm-performance>

**Paper:**

Victoria Caparrós Cabezas and Markus Püschel [Extending the Roofline Model: Bottleneck Analysis with Microarchitectural Constraints](#) Proc. IISWC, pp. 222-231, 2014



# ERM: DAG-Based Performance Model

Application	Model	Parameters (6)																		
<pre> int main(){     int n = 1000;     A=(double*)_mm_malloc(...);     ...     for(int i=0; i&lt;n; i++)         for(int j=0; j&lt;n; j++)             for(int k=0; k&lt;n; k++)                 C[i][j]=A[i][k]*B[k][j]; } </pre>		<table border="1"> <thead> <tr> <th>Node type (<math>x</math>)</th> <th><math>N</math> # nodes</th> </tr> </thead> <tbody> <tr> <td>Computation</td> <td><math>W</math></td> </tr> <tr> <td>A</td> <td><math>W_A</math></td> </tr> <tr> <td>M</td> <td><math>W_M</math></td> </tr> <tr> <td>Memory</td> <td><math>Q</math></td> </tr> <tr> <td>L1</td> <td><math>Q_{L1}</math></td> </tr> <tr> <td>L2</td> <td><math>Q_{L2}</math></td> </tr> <tr> <td>L3</td> <td><math>Q_{L3}</math></td> </tr> <tr> <td>mem</td> <td><math>Q_{mem}</math></td> </tr> </tbody> </table>	Node type ( $x$ )	$N$ # nodes	Computation	$W$	A	$W_A$	M	$W_M$	Memory	$Q$	L1	$Q_{L1}$	L2	$Q_{L2}$	L3	$Q_{L3}$	mem	$Q_{mem}$
Node type ( $x$ )	$N$ # nodes																			
Computation	$W$																			
A	$W_A$																			
M	$W_M$																			
Memory	$Q$																			
L1	$Q_{L1}$																			
L2	$Q_{L2}$																			
L3	$Q_{L3}$																			
mem	$Q_{mem}$																			
Processor	Model	Parameters (22)																		
		<table border="1"> <thead> <tr> <th>Microarchitectural Parameters</th> <th>Unit</th> </tr> </thead> <tbody> <tr> <td>Throughput</td> <td>flops/cycle doubles/cycle instr./cycle</td> </tr> <tr> <td>Latency</td> <td>cycles cycles</td> </tr> <tr> <td>Capacity</td> <td>bytes bytes # slots</td> </tr> </tbody> </table>	Microarchitectural Parameters	Unit	Throughput	flops/cycle doubles/cycle instr./cycle	Latency	cycles cycles	Capacity	bytes bytes # slots										
Microarchitectural Parameters	Unit																			
Throughput	flops/cycle doubles/cycle instr./cycle																			
Latency	cycles cycles																			
Capacity	bytes bytes # slots																			

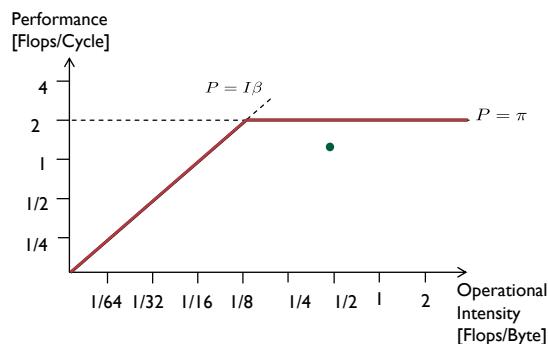
29

## Simple Example: Overlap

### Perfect overlap

$$T = \max(T_{comp}, T_{mem})$$

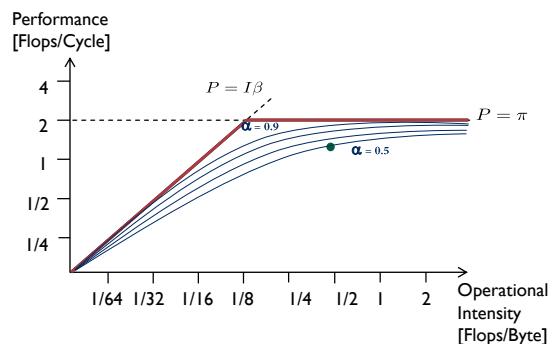
$$P = \min(\beta I, \pi)$$



### Overlap $0 \leq \alpha \leq 1$

$$T = T_{comp} + T_{mem} - \alpha \min(T_{comp}, T_{mem})$$

$$P = \frac{I}{\frac{I}{\pi} + \frac{1}{\beta} - \alpha \min(\frac{I}{\pi}, \frac{1}{\beta})}$$



Challenge 2: How to get overlap information?