

Modeling Stencil Computations on Modern HPC Architectures

Raúl de la Cruz¹ and Mauricio Araya-Polo²

`delacruz@bsc.es`, `mauricio.araya@shell.com`

Barcelona Supercomputing Center (BSC-CNS), Barcelona (Spain)¹
Shell Intl. E&P Inc., Houston, Texas (USA)²

Performance Modeling, Benchmarking and Simulation
of High Performance Computer Systems (SC14)
New Orleans, LA, USA, November 16th, 2014

- 1 Stencil in a Nutshell
- 2 Motivation: Modeling Stencil Performance
- 3 Stencil Performance Model: Base
- 4 Stencil Performance Model: Extended
 - Cache Interference Phenomena
 - Multi-core Considerations
 - Prefetching Approach
 - Stencil Optimizations
- 5 Model Validation
 - Cache Miss Prediction
 - Optimization Techniques
 - Core Efficiency in SMT Mode
- 6 Conclusions & Future Work

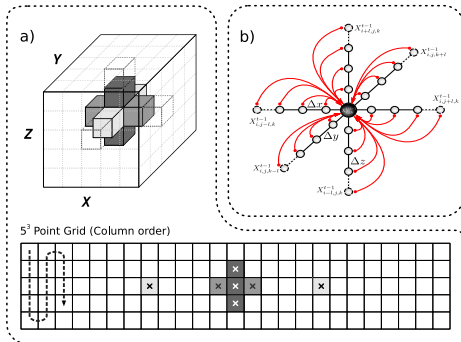
Finite Difference Method

- 1: Domain decomposition of mesh
- 2: **for** $time = 0$ to $time_{end}$ **do**
- 3: Read Input
- 4: Pre-processing
- 5: Inject source
- 6: Apply boundary conditions
- 7: **for** all points in my domain **do**
- 8: Stencil computation (X^t)
- 9: **end for**
- 10: Exchange overlapped points
- 11: Post-processing
- 12: Write Output
- 13: **end for**

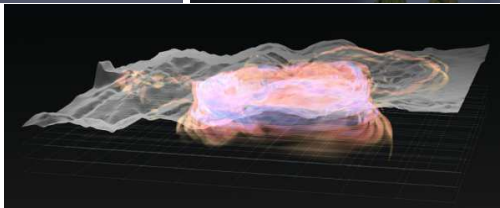
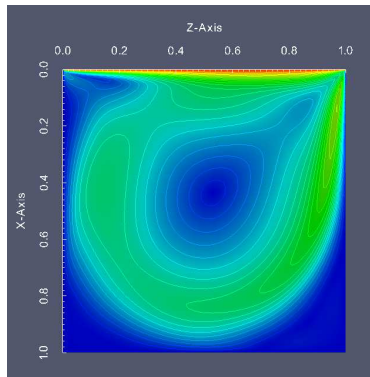
- Load balancing
- Kernel computation
- Intra/inter-node communication

Stencil Computation

```
for  $k = \ell$  to  $Y - \ell$  do
  for  $j = \ell$  to  $X - \ell$  do
    for  $i = \ell$  to  $Z - \ell$  do
       $x_{i,j,k}^t = C_0 * x_{i,j,k}^{t-1}$ 
        +  $C_{Z1} * (x_{i-1,j,k}^{t-1} + x_{i+1,j,k}^{t-1}) + \dots + C_{Z\ell} * (x_{i-\ell,j,k}^{t-1} + x_{i+\ell,j,k}^{t-1})$ 
        +  $C_{X1} * (x_{i,j-1,k}^{t-1} + x_{i,j+1,k}^{t-1}) + \dots + C_{X\ell} * (x_{i,j-\ell,k}^{t-1} + x_{i,j+\ell,k}^{t-1})$ 
        +  $C_{Y1} * (x_{i,j,k-1}^{t-1} + x_{i,j,k+1}^{t-1}) + \dots + C_{Y\ell} * (x_{i,j,k-\ell}^{t-1} + x_{i,j,k+\ell}^{t-1})$ 
    end for
  end for
end for
```



Stencil in the Academia & Industry



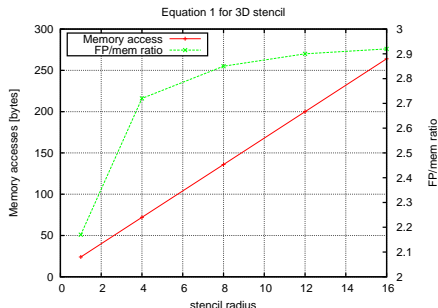
Stencil Performance Challenges

Two main challenges:

1 Low FLOPs/Memory ratio

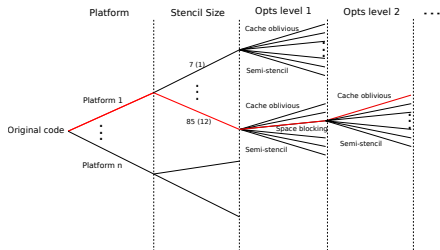
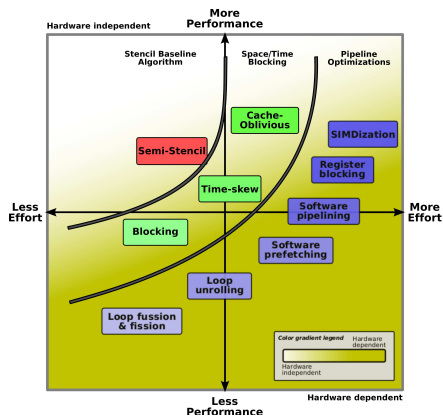
$$\begin{aligned} FP/Cache_{Classical} &= \frac{\text{FloatingPoint Operations}}{\text{Data Cache Accesses}} \\ &= \frac{2 * \text{MultiplyAdd Instructions}}{X^{t-1} \text{ Loads} + X^t \text{ Stores}} \\ &= \frac{2 * 2 * dim * \ell + 1}{\underbrace{(2 * (dim - 1) * \ell + 1)}_{X^{t-1} \text{ Loads}} + \underbrace{(1)}_{X^t \text{ Stores}}} \\ &= \frac{4 * dim * \ell + 1}{2 * dim * \ell - 2 * \ell + 2} \end{aligned} \quad (1)$$

2 Low data reuse



- 1 Stencil in a Nutshell
- 2 Motivation: Modeling Stencil Performance
- 3 Stencil Performance Model: Base
- 4 Stencil Performance Model: Extended
 - Cache Interference Phenomena
 - Multi-core Considerations
 - Prefetching Approach
 - Stencil Optimizations
- 5 Model Validation
 - Cache Miss Prediction
 - Optimization Techniques
 - Core Efficiency in SMT Mode
- 6 Conclusions & Future Work

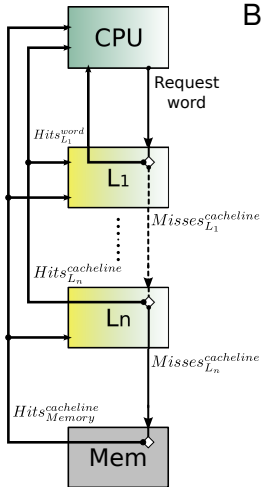
Motivation: Modeling Stencil Performance



- Which one/combination of these should be implemented?
- How much performance improvement can be expected?
- Brute force (autotuning) is too expensive (combinatorial explosion, code management) then modeling seems to be a valid alternative

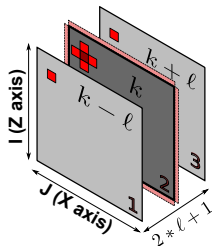
- 1 Stencil in a Nutshell
- 2 Motivation: Modeling Stencil Performance
- 3 Stencil Performance Model: Base**
- 4 Stencil Performance Model: Extended
 - Cache Interference Phenomena
 - Multi-core Considerations
 - Prefetching Approach
 - Stencil Optimizations
- 5 Model Validation
 - Cache Miss Prediction
 - Optimization Techniques
 - Core Efficiency in SMT Mode
- 6 Conclusions & Future Work

Stencil Performance Model: Base Model I



Basic considerations:

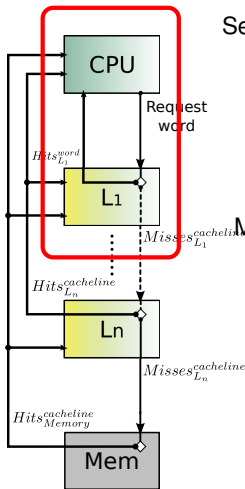
- Computation bottleneck is negligible (low FP/Byte) and reads dominate
- No cache interferences between instructions and data
- P_{read} ($2 * \ell + 1$) Z-X planes from \mathcal{X}^{t-1} to compute one \mathcal{X}^t output plane (P_{write})
- $S_{total} = P_{read} \times S_{read} + P_{write} \times S_{write}$, being $S_{read} = I \times J$ and $S_{write} = I \times J$
- Three memory groups are established



$$T_{total} = \underbrace{T_{L1}}_{first} + \underbrace{\dots + T_{Li} + \dots + T_{Ln}}_{intermediate} + \underbrace{T_{Memory}}_{last}$$

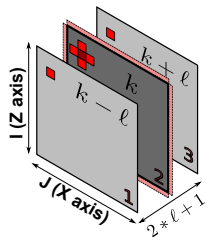
What are the $Hits_{Li}$ and $Misses_{Li}$ for each hierarchy level?

Stencil Performance Model: Base Model II



Several data are involved in a stencil comput.:

- Grid points (Input: \mathcal{X}^{t-1} , Output: \mathcal{X}^t)
- Contribution weights ($C_{Z,X,Y,0}$)
- Indices to access grid points (i, j, k)



Misses depend on II , JJ sizes and ℓ parameter:

$$Misses_{Li}^{cline} = \lceil II/W \rceil * JJ * KK * nplanes_{Li}$$

$$Hits_{Li}^{cline} = Misses_{Li-1}^{cline} - Misses_{Li}^{cline}$$

$$T_{Li}^{cline} = cacheline / Bw_{Li}^{read}$$

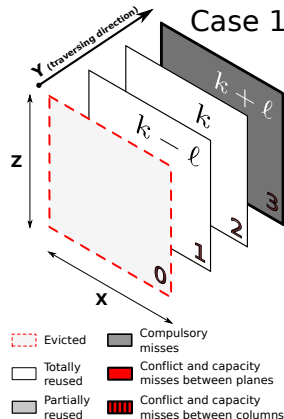
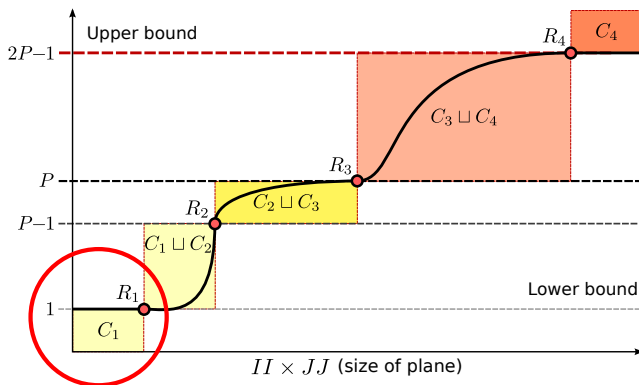
$$T_{Li} = Hits_{Li}^{cline} * T_{Li}^{cline}$$

- $nplanes_{Li}$: $II \times JJ$ planes read from $Li + 1$ for each k iteration

(Basic model description in ICCS-iWAPT2011 paper)

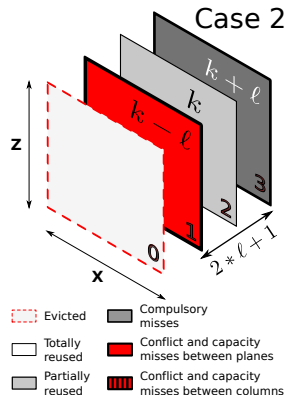
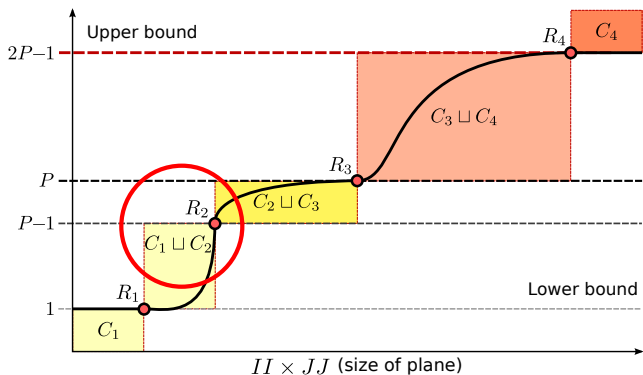
Stencil Model: Read Miss Cases & Rules

$$nplanes_{Li}(II, JJ) = \begin{cases} C_1 : 1, & \text{if } R_1 \\ C_1 \sqcup C_2 : (1, P_{read} - 1], & \text{if } \neg R_1 \wedge R_2 \\ C_2 \sqcup C_3 : (P_{read} - 1, P_{read}], & \text{if } \neg R_2 \wedge R_3 \\ C_3 \sqcup C_4 : (P_{read}, 2P_{read} - 1], & \text{if } \neg R_3 \wedge \neg R_4 \\ C_4 : 2P_{read} - 1, & \text{if } R_4, \end{cases}$$



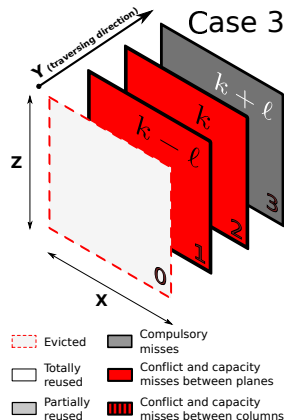
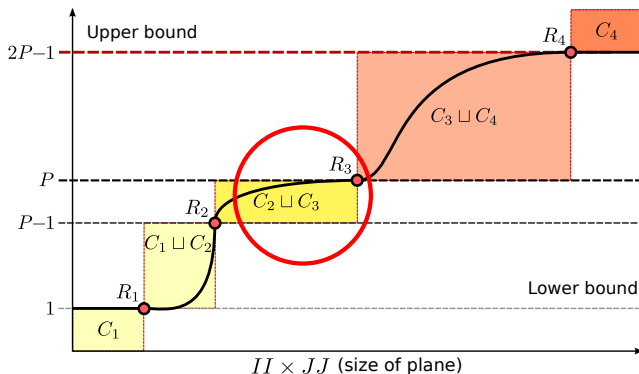
Stencil Model: Read Miss Cases & Rules

$$nplanes_{Li}(II, JJ) = \begin{cases} C_1 : 1, & \text{if } R_1 \\ C_1 \sqcup C_2 : (1, P_{read} - 1], & \text{if } \neg R_1 \wedge R_2 \\ C_2 \sqcup C_3 : (P_{read} - 1, P_{read}], & \text{if } \neg R_2 \wedge R_3 \\ C_3 \sqcup C_4 : (P_{read}, 2P_{read} - 1], & \text{if } \neg R_3 \wedge \neg R_4 \\ C_4 : 2P_{read} - 1, & \text{if } R_4, \end{cases}$$



Stencil Model: Read Miss Cases & Rules

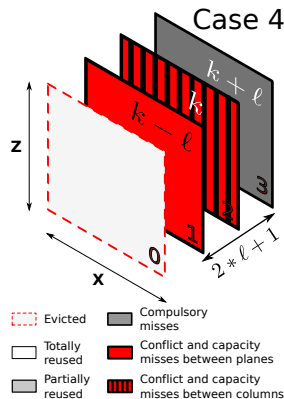
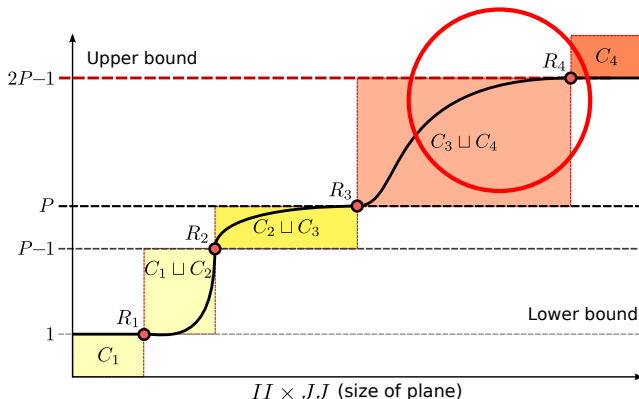
$$nplanes_{Li}(II, JJ) = \begin{cases} C_1 : 1, & \text{if } R_1 \\ C_1 \sqcup C_2 : (1, P_{read} - 1], & \text{if } \neg R_1 \wedge R_2 \\ C_2 \sqcup C_3 : (P_{read} - 1, P_{read}], & \text{if } \neg R_2 \wedge R_3 \\ C_3 \sqcup C_4 : (P_{read}, 2P_{read} - 1], & \text{if } \neg R_3 \wedge \neg R_4 \\ C_4 : 2P_{read} - 1, & \text{if } R_4, \end{cases}$$



- Evicted
- Compulsory misses
- Totally reused
- Conflict and capacity misses between planes
- Partially reused
- Conflict and capacity misses between columns

Stencil Model: Read Miss Cases & Rules

$$nplanes_{Li}(II, JJ) = \begin{cases} C_1 : 1, & \text{if } R_1 \\ C_1 \sqcup C_2 : (1, P_{read} - 1], & \text{if } \neg R_1 \wedge R_2 \\ C_2 \sqcup C_3 : (P_{read} - 1, P_{read}], & \text{if } \neg R_2 \wedge R_3 \\ C_3 \sqcup C_4 : (P_{read}, 2P_{read} - 1], & \text{if } \neg R_3 \wedge \neg R_4 \\ C_4 : 2P_{read} - 1, & \text{if } R_4, \end{cases}$$

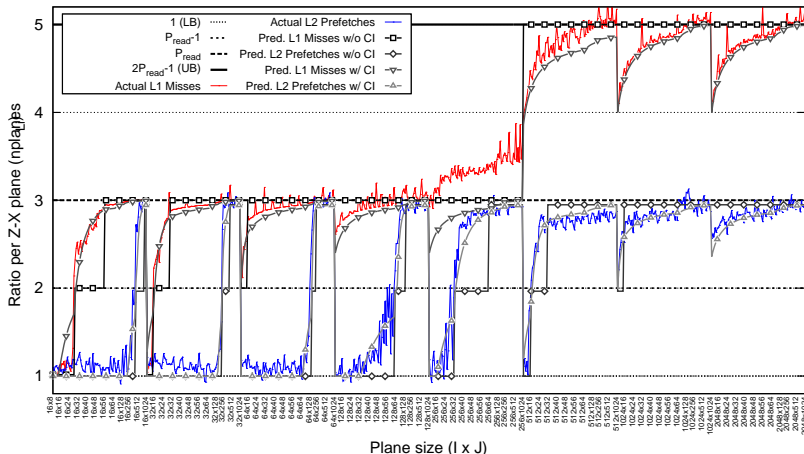


- 1 Stencil in a Nutshell
- 2 Motivation: Modeling Stencil Performance
- 3 Stencil Performance Model: Base
- 4 Stencil Performance Model: Extended**
 - Cache Interference Phenomena
 - Multi-core Considerations
 - Prefetching Approach
 - Stencil Optimizations
- 5 Model Validation
 - Cache Miss Prediction
 - Optimization Techniques
 - Core Efficiency in SMT Mode
- 6 Conclusions & Future Work

Stencil Model: Cache Interference Phenomena

- Add full 3C (compulsory, conflict and capacity) misses detection
- Convert the model from a discrete to a continuum space

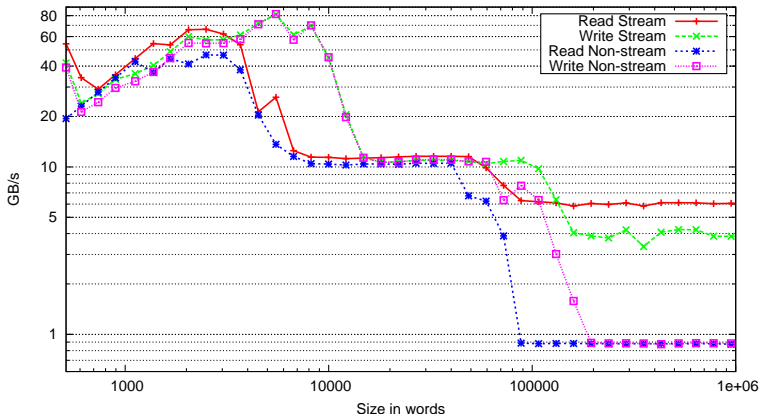
$$P(i) = \frac{\overbrace{II \times JJ}^{\text{whole plane}} - \overbrace{II \times (P_{\text{read}} - 1)}^{\text{columns reuse}}}{II \times JJ} \in [0, 1] \rightarrow nplanes_{Li}' = nplanes_{Li} \times P(i)$$



Stencil Model: from Single-core to Multi-core

STREAM2 benchmark must mimic stencil environment:

- Read Bw (DOT kernel), Write Bw (FILL kernel)
- Number of threads and their pinning to cores
- Memory alignment, temporal or non-temporal writes, SISD/SIMD ISA
- Shared resources $\rightarrow size_{Li} = N/nthreads_{core}$

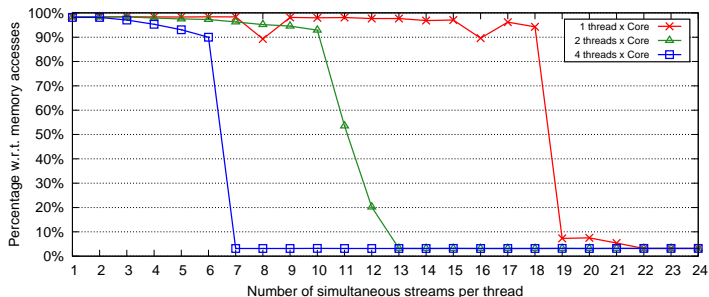


Stencil Model: Prefetching Effect

 Gabriel Marin, Collin McCurdy, and Jeffrey S. Vetter.

Diagnosis and optimization of application prefetching performance. (ICS '13)

$$DC_{effectiveness} = DC_Req_PF / DC_Req_All \in [0, 1]$$



$$nplanes_{Li}^S = nplanes_{Li} \times DC_{effectiveness}$$

$$nplanes_{Li}^{NS} = nplanes_{Li} \times (1 - DC_{effectiveness})$$

Stencil Model: Spatial Blocking Optimization

Traverse the domain in $TI \times TJ \times TK$ blocks:

$$\begin{aligned} I &= \lceil TI/W \rceil \times W, & J &= TJ, & K &= TK, \\ II &= \lceil (TI + 2 \times \ell)/W \rceil \times W, & JJ &= TJ + 2 \times \ell, & KK &= TK + 2 \times \ell. \end{aligned}$$

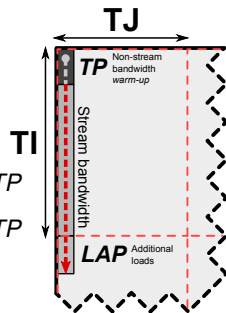
New II and JJ are used for $R_{1,2,3,4}$ to estimate $nplanes_{Li}$:

$$Misses_{Li}^{[S,NS]} = \lceil II/W \rceil \times JJ \times KK \times nplanes_{Li}^{[S,NS]} \times NB,$$

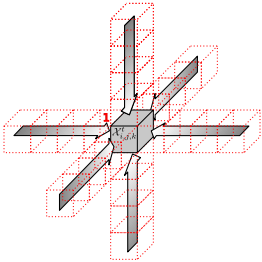
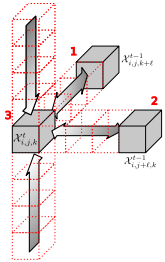
Consider penalties on caches with prefetching (may disrupt memory access and increase data transfers):

$$Misses_{Li}^{NS} \pm TP \times JJ \times KK \times nplanes_{Li}^{NS} \times NB, \quad \text{if } II/W \geq TP$$

$$Misses_{Li}^S \pm LAP \times JJ \times KK \times nplanes_{Li}^S \times NB, \quad \text{if } II/W \geq TP$$



Stencil Model: Semi-stencil Optimization

| Dataset accesses | Classical Stencil | Partial <i>Semi-stencil</i> (X and Y axis) |
|---------------------------|---|--|
| Structure of the stencil |  |  |
| \mathcal{X}^{t-1} loads | $2 * (dim - 1) * \ell + 1$ | $(dim - 1) * \ell + 1$ |
| \mathcal{X}^t loads | 0 | $dim - 1$ |
| \mathcal{X}^t stores | 1 | dim |

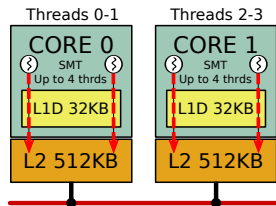
Semi-stencil can be modeled by setting the Z-X planes read/written

$$\begin{aligned}
 P_{read} &= \overbrace{\ell + 1}^{\mathcal{X}^{t-1}} + \overbrace{1}^{\mathcal{X}^t}, & P_{write} &= \overbrace{2}^{\mathcal{X}^t}, & \text{if } \neg R_4 \\
 P_{read} &= \ell + 1 + 2, & P_{write} &= 3, & \text{if } R_4
 \end{aligned}$$

- 1 Stencil in a Nutshell
- 2 Motivation: Modeling Stencil Performance
- 3 Stencil Performance Model: Base
- 4 Stencil Performance Model: Extended
 - Cache Interference Phenomena
 - Multi-core Considerations
 - Prefetching Approach
 - Stencil Optimizations
- 5 Model Validation**
 - Cache Miss Prediction**
 - Optimization Techniques**
 - Core Efficiency in SMT Mode**
- 6 Conclusions & Future Work

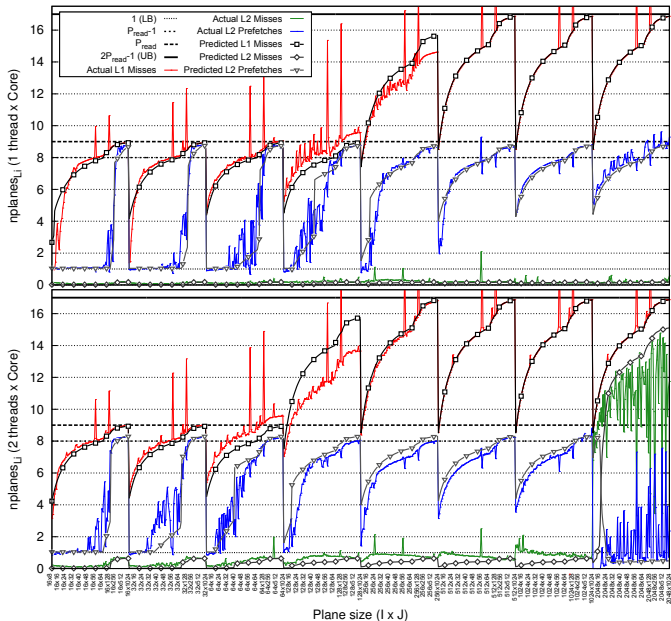
Model Validation: Experiments on Intel Xeon Phi

| Parameters | Range of values |
|--|--|
| Naive sizes ($I \times J \times K$) | $8 \times 8 \times 128 \dots 2048 \times 1024 \times 128$ |
| Rivera sizes ($I \times J \times K$) | $512 \times 2048 \times 128$ |
| Stencil sizes (ℓ) | 1, 2, 4 and 7 (7, 13, 25 and 43-point) |
| Algorithms | {Naive, Rivera} \times {Classical, <i>Semi-stencil</i> } |
| Block sizes (T_I and T_J) | {8, 16, 24, 32, 64, 128, 256, 512, 1024, 1536, 2048} |
| SMT configuration | 1, 2 and 4 threads x Core |

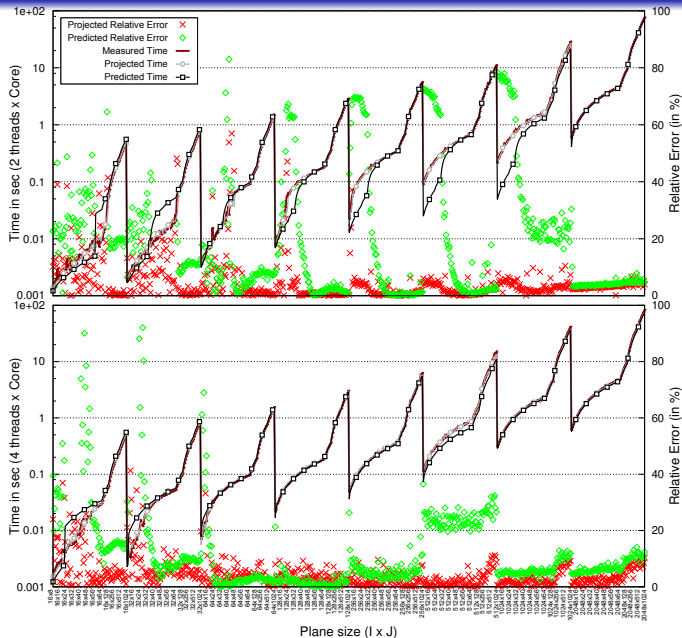


| Description | Intel Xeon Phi Events | Time Cost Formulas |
|---------------|-----------------------|--|
| Cycles | CPU_CLK_UNHALTED | $T_{L1} = (L1 \text{ Reads} - L1 \text{ Misses}) \times Bw_{L1}^{cline}$ |
| L1 Reads | VPU_DATA_READ | $T_{L2} = Bw_{L2}^{cline} \times (L1 \text{ Misses} - L2 \text{ Misses} - (L2 \text{ Prefetches} - L2 \text{ Writes} \times \text{Pref Eff}))$ |
| L1 Misses | VPU_DATA_READ_MISS | |
| L2 Misses | L2_DATA_READ_MISS_ | $T_{Mem} = L2 \text{ Misses} \times Bw_{Mem}^{NS} + Bw_{Mem}^S \times (L2 \text{ Prefetches} - L2 \text{ Writes} \times \text{Pref Eff})$ |
| | MEM_FILL | |
| L2 Prefetches | HWP_L2MISS | $T_{Writes} = L2 \text{ Writes} \times \text{Pref Eff} \times Bw_{Write}^S$ |
| L2 Writes | L2_WRITE_HIT | $+ L2 \text{ Writes} \times (1 - \text{Pref Eff}) \times Bw_{Write}^{NS}$ |

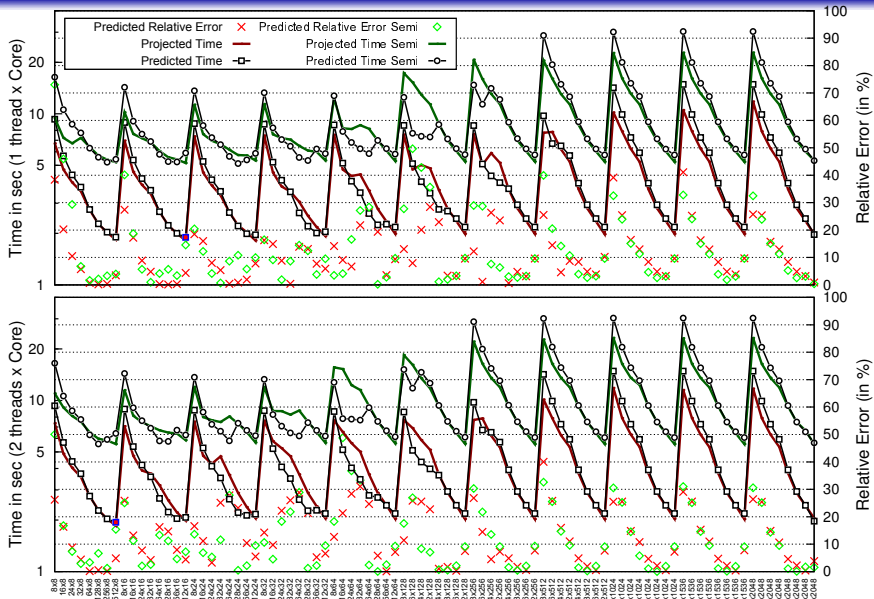
Model Validation: Cache Misses (Naive case, $\ell = 4$)



Model Validation: Time Prediction (Naive case, $\ell = 7$)

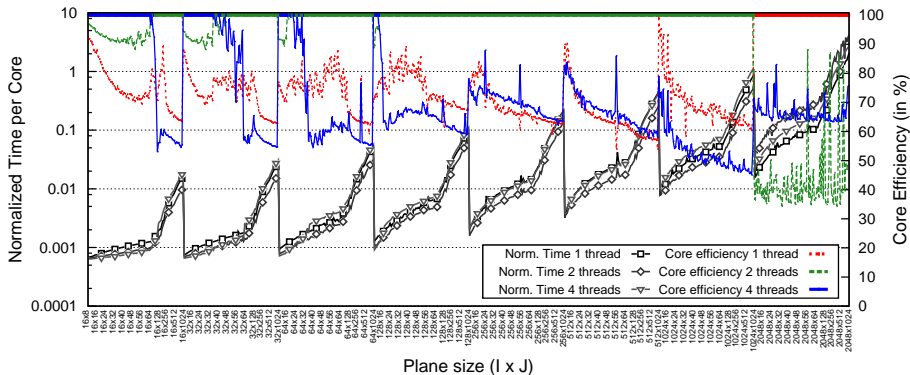


Model Validation: Spatial Blocking & Semi-stencil ($\ell = 1, \ell = 7$)



Model Validation: Core Efficiency in SMT Mode

$$Core_{efficiency}^{SMT_i} = \frac{\min(\tau^{SMT_1}, \dots, \tau^{SMT_n})}{\tau^{SMT_i}} \in [0, 1]$$



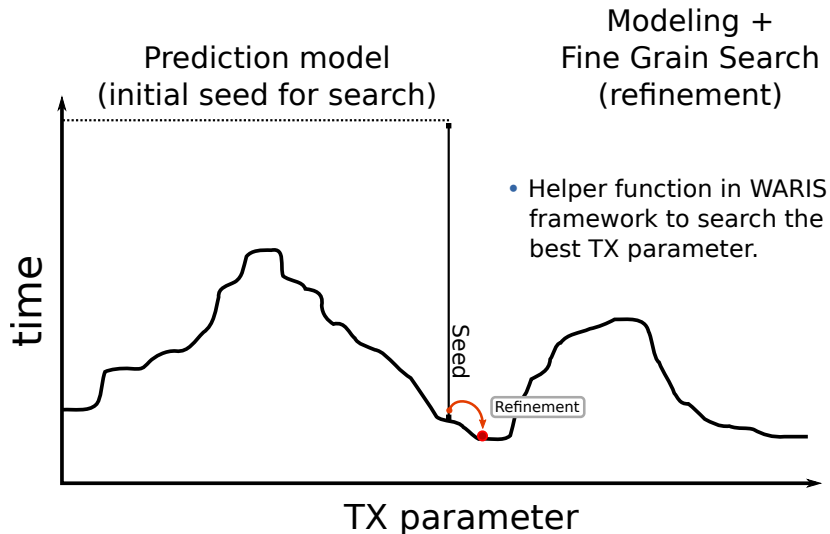
- This model includes multi- and many-core support
- Robust hardware prefetching modeling (*prefetching effectiveness*)
- Considers cache interference phenomena (3C misses)
- Extended with stencil optimizations (blocking and Semi)
- Most results have a high accuracy (rel. error 5-15%)
- Core efficiency predictor (best SMT configuration)

Thank you for your attention!

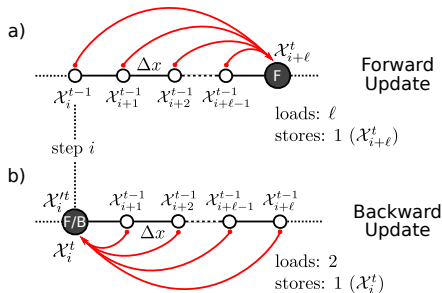


Copyright 2013. Barcelona Supercomputing Center - BSC

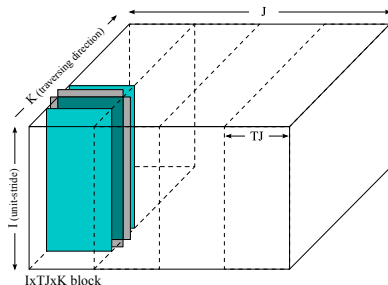
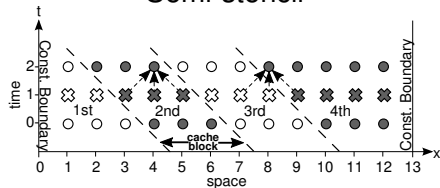
Auto-tuning blocking parameters



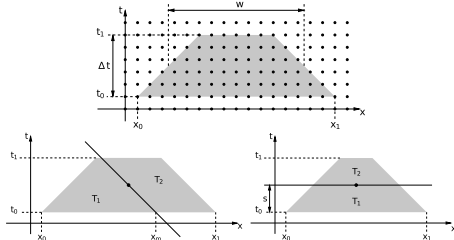
State of the Art - The most popular



Semi-stencil



Rivera



Cache-oblivious

Stencil Model: Read Miss Cases & Rules

