# SMAT：
# An Input Adaptive Auto-Tuner for Sparse Matrix-Vector Multiplication

**Jiajia Li, Guangming Tan, Mingyu Chen, Ninghui Sun**

Institute of Computing Technology,
Chinese Academy of Sciences

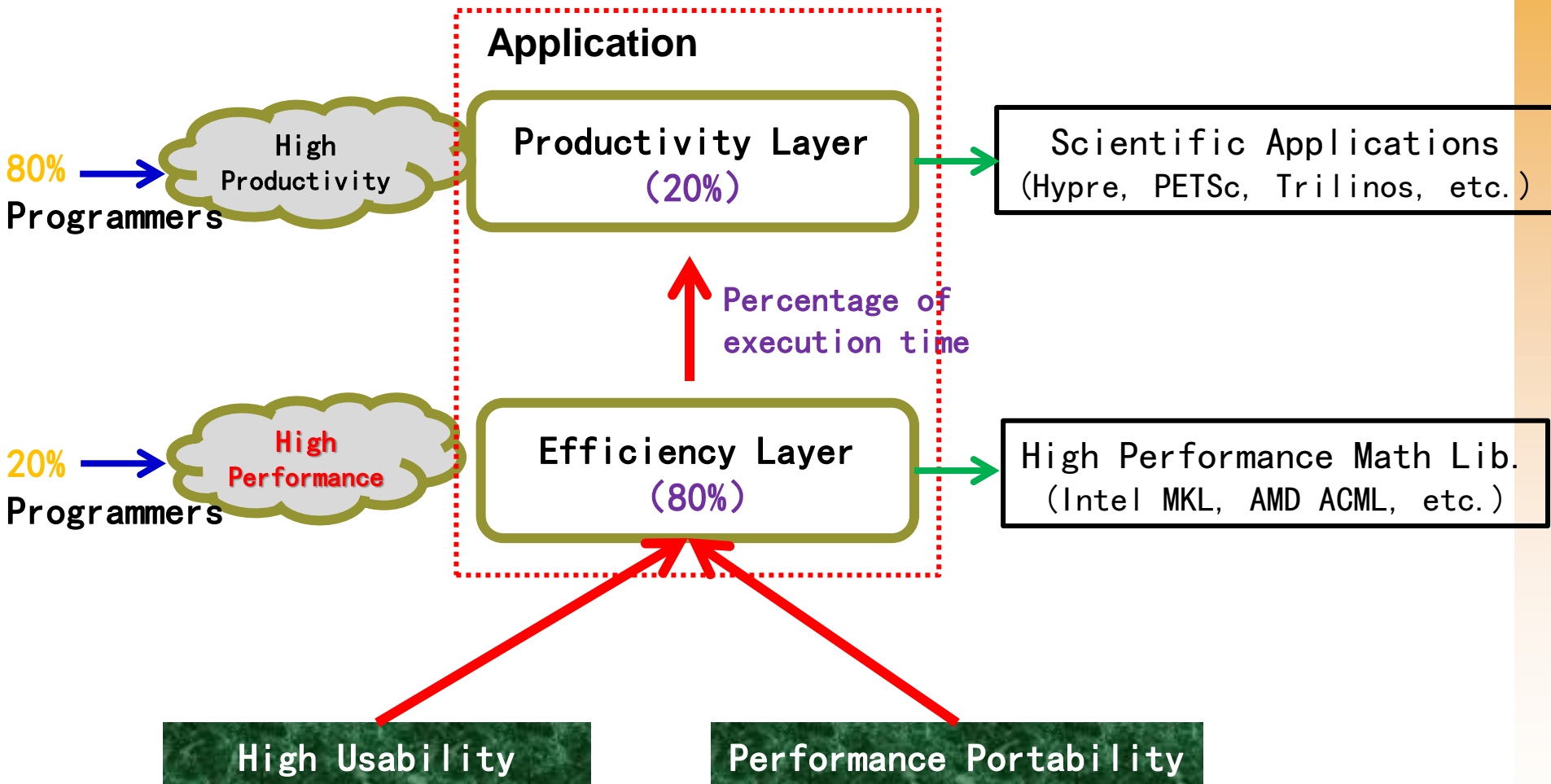Institute of Computing Technology,Chinese Academy of Sciences

# Contents

# High Performance Computational Software Development-1

**Computational Science**

**Data Science**

ITER

Climate

Oil

**Exascale**

**Big Data**

Social network

National Security

System Biology

## Sparse Linear System

Sparse Matrices

| Protein | FEM / Spheres | FEM / Cantilever | Wind Tunnel | FEM / Harbor | QCD |
| FEM / Ship | Economics | Epidemiology | FEM / Accelerator | Circuit | webbase |

# Sparse Matrix

◆ **Diverse Application Background**

◆ **Different Solving Methods**

$\Rightarrow$

**Diff. Nonzero Distribution Structure**

$\Downarrow$

● **Kinds of Sparse Matrices**

– Diagonal Matrix

– "Slim" Matrix

– "Fat" Matrix

– Power-Law Matrix

– Matrix with Dense Blocks

– …



2

Diagonal
"pcrystk02"

"Slim"
"Bfly"

"Fat"
"crankseg_2"

Power-Law
"roadNet-CA"

222

number of nonzeros per row

## Storage Formats

$$A = \begin{bmatrix} 1 & 5 & 0 & 0 \\ 0 & 2 & 6 & 0 \end{bmatrix}$$

SpMV: solve $Y = AX + Y$, where **A** is a sparse matrix, **X** and **Y** are dense vectors.

data [1 5 2 6 8 3 7 9 4]

```
sum = x[indices[jj]] * data[jj];
y[i] = sum;
}
```
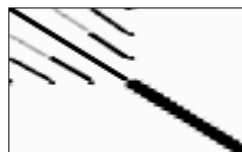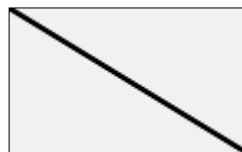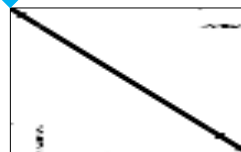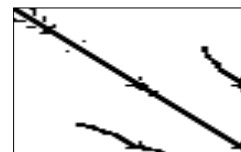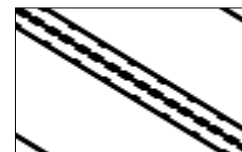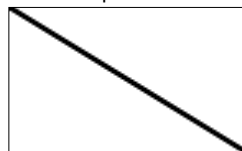
**(a) CSR SpMV**

row [0 0 1 1 2 2 2 3 3]
col [0 1 1 2 0 2 3 1 3]
data [1 5 2 6 8 3 7 9 4]

```
for (i = 0; i <num_nonzeros; i++)
{
  y[rows[i]] = data[i] * x[cols[i]];
}
```

**(b) COO SpMV**

offsets [-2 0 1]

$$data = \begin{bmatrix} * & 1 & 5 \\ * & 2 & 6 \\ 8 & 3 & 7 \\ 9 & 4 & * \end{bmatrix}$$

```
for( i = 0; i <num_diags; i++){
  k = offsets[i];   //diagonal offset
  Istart = max((0,-k);
  Jstart = max(0, k);
  N = min(num_rows - Istart, num_cols -
Jstart);
  for( n = 0; n < N; n++){
    y_[Istart+n] =data[Istart+i*stride+ n] *
    x[Jstart + n];
  }
}
```

**(c) DIA SpMV**

$$Indices = \begin{bmatrix} 0 & 1 & * \\ 1 & 2 & * \\ 0 & 1 & 2 \\ 1 & 3 & * \end{bmatrix}$$

$$data = \begin{bmatrix} 1 & 5 & * \\ 2 & 6 & * \\ 8 & 3 & 7 \\ 9 & 4 & * \end{bmatrix}$$

```
for(n = 0; n <max_ncols; n++)
{
  for(i = 0; i <num_rows; i++)
  y[i] =
data[n*num_rows+i] *
      x[indices[n*num_rows+i]];
}
```

**(d) ELL SpMV**

# Motivation

**Sparse solvers mainly use ONE storage format CSR**

  ■ Hypre (LLNL), PETSc (ANL), Trilinos (SNL)

**Low Performance**

**GAP!**

(*csr* is the only one exposed to users)

smat_**<precision > <csr><operation>** ( )

mkl_**<precision > <format> <operation>** ( )

(variants of *format* )

**High Performance
High Productivity**

**Libraries provide complicated interfaces to users**

  ■ MKL (Intel), OSKI (UCB), SpBLAS (UTK)

**Low Productivity**

Observation 1: The optimal formats are diverse for sparse matrices from different application areas.

Observation 2: Different formats are needed in different stages of one application during runtime.



## Algebraic Multigrid (AMG) Solver



# Performance Gap: 10x!

8

◆ **Offline**

- ■ Extract application
- ■ Determine feature
- ■ Build feature databas (including the optima implementation)
- ■ Summarize the rules
- ■ Choose the optimal in architecture characte

◆ **Online**

- ■ Extract parameter values of the input matrix
- ■ Predict the optimal algo. & impl. based on model

**Sparse App.**
- Matrix Dimension
- Diagonal Situation
- Nonzero Distribution
- Nonzero Fill Ratio

**Graph App.**
- Dimension
- Degree
- Distribution
- Diameter
- Un-/directed
- Power-Law
- Connectivity

Offline → Online ·····> Call

Application → Architecture

Extract App. Feature → Extract Arch. Feature

Parameter Set → Various Implementations

Get Para. Value (Learning Set) → Choose

Feature Database → Optimal Implementation

Summarize Features

Model → Predict → Optimal Algo. & Impl.

Online | Extract Features ← Application

Algorithm / Data Structure / Architecture

# SMAT Framework

## SMAT

**Example**
- TLB size
- Cache size
- Reg. Size
- Prefecth
- SIMDization
- Branch
- Multi-thread

**Example**
- Matrix Dim.
- Diagonal Situation
- Nonzero Ratio
- Nonzero Dist.
- Power-Law

Offline → **Offline**
Online → **Online**

Search

**Arch. Parm.**

**Kernel Library**

**Optimal Kernel**

**Exec. & Measure**

Link

< TH

Learning Model

Optimal SpMV

**Kernel Opt.**

**Spl Behavior**

**App. Parm.**

**Feature Database**

Determine Parm. Value

Data Mining (C5.0)

**Feature Extraction**

SMAT_x**CSR**_SpMV

Offline

Online

**Intel MKL**

**SMAT**

mkl_xcsrgemv
mkl_xdiagemv
mkl_xbsrgemv
mkl_xcscmv
mkl_xcoogemv
mkl_xskymv

→ SMAT_xCSR_SpMV

IT'S Easy

So Easy!

**Generally:**
- ◆ Diagonal Matrix ↔ DIA Format
- ◆ "Slim" Matrix ↔ ELL Format
- ◆ "Fat" Matrix ↔ CSR Format
- ◆ Power-Law Matrix ↔ COO Format

| Parameter | | Meaning | Formula | DIA | ELL | CSR | COO |
|---|---|---|---|---|---|---|---|
| Matrix Dimension | M | the number of rows | - | √ | √ | √ | √ |
| | N | the number of columns | - | √ | √ | √ | √ |
| Diagonal Situation | Ndiags | the number of diagonals | - | ↓ | | | |
| | NTdiags_ratio | the ratio of "true" diagonals to total diagonals | $NTdiags\_ratio = \frac{number\ of\ "true\ diagonals"}{Ndiags}$ | ↑ | | | |
| Nonzero Distribution | NNZ | the number of nonzeros | - | √ | √ | √ | √ |
| | aver_RD | the number of nonzeros per row | $aver\_RD = \frac{NNZ}{M}$ | √ | √ | √ | √ |
| | max_RD | the maximum number of nonzeros per row | $max\_RD = \max_1^M \{number\ of\ nonzeros\ per\ row\}$ | | ↓ | | |
| | var_RD | the variation of the number of nonzeros per row | $var\_RD = \frac{\Sigma_1^M |row\_degree - aver\_RD|^2}{M}$ | | | | ↓ |
| Nonzero Ratio | ER_DIA | the ratio of nonzeros in DIA data structure | $ER\_DIA = \frac{NNZ}{Ndiags \times M}$ | ↑ | | | |
| | ER_ELL | the ratio of nonzeros in ELL data structure | $ER\_ELL = \frac{NNZ}{max\_RD \times M}$ | | ↑ | | |
| Power-Law | R | a factor of power-law distribution | $P(k) \sim k^{-R}$ | | | | [1, 4] |

"√" shows the parameter is useful for all formats.
"↑/↓" indicates the larger (smaller) the parameter value is, the format shows more benefit.

# Feature Extraction Process

◆ **Divide the matrix set based on the optimal storage format, and measure SpMV performance on them**

◆ **Draw the value distribution graph to each parameter, and find the regulations.**



Perf. Beneficial values

"DIA_best" Matrix Set — DIA-"Ndiags"

Ndiags with small value

DIA format earns most benefit

Ndiags values

[1,4]

(e) R

| Matrix Partition | DIA_best | ELL_best | CSR_best | COO_best |
|---|---|---|---|---|
| Size of Sub-set | 206 | 169 | 1496 | 507 |

## ◆ Scoreboard Strategy

- Choose typical matrix set to test each algo. and impl.
- Record the performance value on scoreboard
- The optimal impl. for each algo. are recorded in index table

**Example**
- TLB size
- cache size
- register size
- prefetch
- SIMDization
- branch
- threading policy

**Example**
- matrix dimension
- diagonal situation
- nonzero ratio
- nonzero distribution
- power-law

One record

{ M, N, Ndiags, Ntdiags_ratio, NNZ, max_RD, var_RD, ER_DIA, ER_ELL, R, Best_Format }

◆**Belong to data mining problems**

- Classification problem
- Target Attribute: "Best_format"

$$f\left(\overrightarrow{x_1}, \overrightarrow{x_2}, \cdots, \overrightarrow{x_n}, \overrightarrow{TH}\right) \rightarrow C_n(DIA, ELL, CSR, COO)$$

$\overrightarrow{x_i}$: value of each record, $\overrightarrow{TH}$: threshold of each parameter, $C_n(DIA, ELL, CSR, COO)$: one of the four formats. `

◆**Build model**

- Use ruleset to represent model
- Add rule confidence

**Data Mining Process**

## Classification Algorithm

### Training Data

| M | N | Ndiags | Ntdiags _ratio | ... | Best _format |
|---|---|--------|----------------|-----|--------------|
| 36476 | 36476 | 1328 | 0.001 | | ELL |
| 23560 | 23560 | 33 | 0.515 | | DIA |
| 12M | 12M | 1.8M | ~0 | | COO |
| ... | ... | ... | ... | | ... |

### Rules

IF Ndiags <= 39
AND ER_DIA > 0.3
AND ER_ELL <= 0.9
THEN
Best_format = DIA

◆ **Platforms：**

- Intel Xeon X5680
- AMD Opteron 6168

◆ **Matrix set: The University of Florida Sparse Matrix Collection**

- Learning Set：2055
- Testing Set：331，represented by 16 matrices

**Representative Matrices**

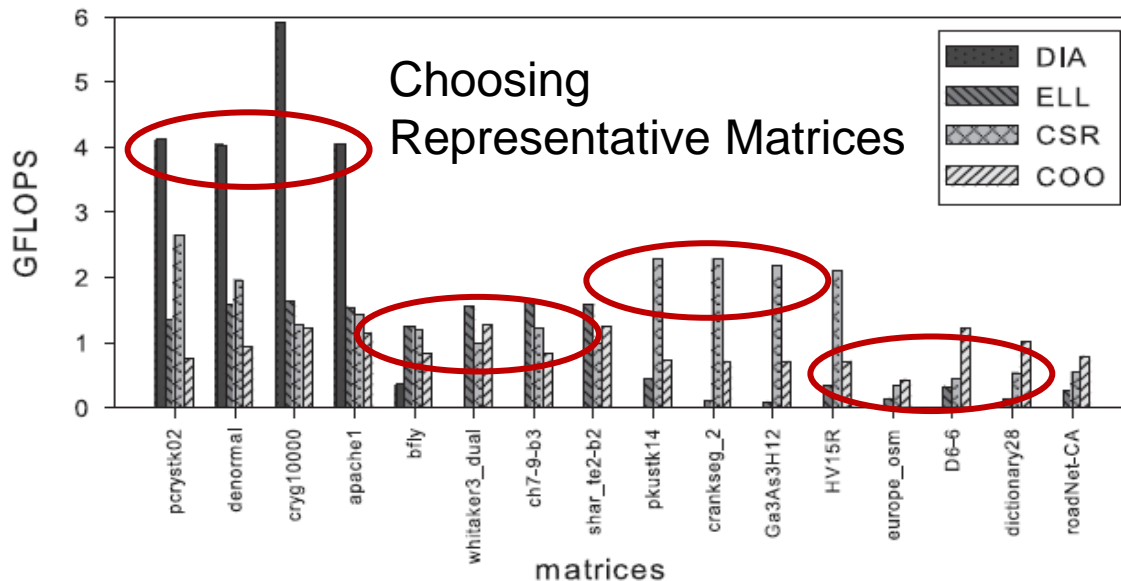| No. | Graph | Name | Dimensions | Nonzeros (NNZ / M) | Application area |
|-----|-------|------|------------|--------------------|------------------|
| 1 | | pcrystk02 | 14K×14K | 491K (35) | duplicate materials problem |
| 2 | | d... | | | ...-example |
| 3 | | cryg10000 | 10K×10K | ...M (5) | materials problem |
| 4 | | apache1 | 81K×81K | 311K (4) | structural problem |
| 5 | | bfly | 49K×49K | 98K (2) | undirected graph sequence |
| 6 | | whita..._dual | | | ...3D problem |
| 7 | | ch7-9-b3 | 106K×18K | 423K (4) | combinatorial problem |
| 8 | | shar_te2-b2 | 200K×17K | 601K (3) | combinatorial problem |
| 9 | | pkustk14 | 152K×152K | 15M (98) | structural problem |
| 10 | | crankseg... | | | ...uctural problem |
| 11 | | Ga3As3H12 | 61K×61K | ...M (97) | theoretical/quantum chemistry |
| 12 | | HV15R | 2M×2M | 283M (140) | computational fluid dynamics |
| 13 | | europe_osm | 51M×51M | 108M (2) | undirected graph |
| 14 | | | | | ...problem |
| 15 | | dictionary28 | 53K×53K | 178K (3) | undirected graph |
| 16 | | roadNet-CA | 2M×2M | 6M (3) | undirected graph |

**Diagonal Mat.**

**"Slim" Mat.**

**"Fat" Mat.**

**Power-Law Mat.**



Choosing Representative Matrices

GFLOPS — DIA, ELL, CSR, COO

matrices: pcrystk02, denormal, cryg10000, apache1, bfly, whitaker3_dual, ch7-9-b3, shar_te2-b2, pkustk14, crankseg_2, Ga3As3H12, HV15R, europe_osm, D6-6, dictionary28, roadNet-CA
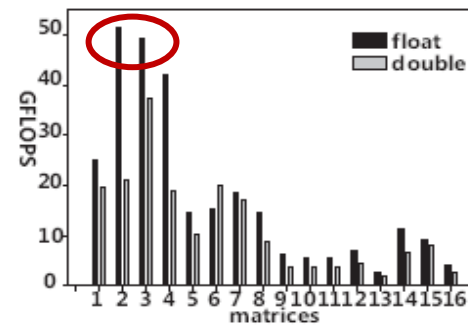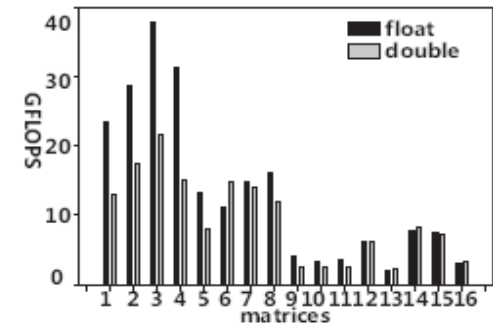
◆ **SMAT Auto-tuning**

◆ **Optimized SpMV kernels**

- Assembling opt.
  - Loop unrolling
  - SIMDization
- Multi-threading on task level
  - Allocate a sub-block to each thread
  - Independently choose the optimal algo. & impl. on each sub-block

Performance
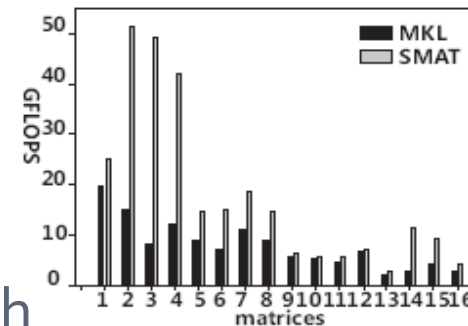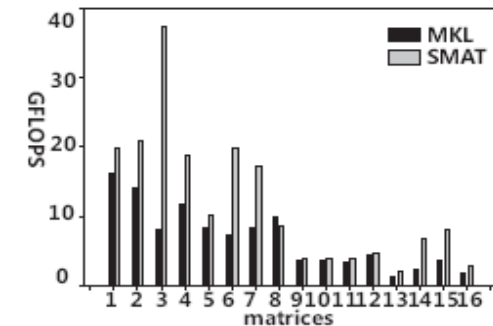


(a) Intel    (b) AMD

Compared with MKL



(a) float    (b) double

~3X Speedup on average

# Analysis on Accuracy and Overhead

◆ **Analyze the prediction procedure and accuracy on 16 representative matrices**

◆ **Overhead**

- When the model predicts right, small overhead (about **2** CSR-SpMV)
- Wrong predict, execute & measure mudule used, the overhead is more than **10** CSR-SpMV   (OSKI:40+; clSpMV: ~10）
- When a matrix is used hundreds of times in an iterative method, the overhead can be overlapped.

| Matrix Number | Matrix Name | Model Prediction Format | Execution | SMAT Prediction Format | Actual Best Format | Model Accuracy | SMAT Overhead (times of CSR-SpMV) |
|---|---|---|---|---|---|---|---|
| 1 | pcrystk02 | DIA | - | DIA | DIA | R | 2.28 |
| 2 | denormal | DIA | - | DIA | DIA | R | 2.09 |
| 3 | cryg10000 | DIA | - | DIA | DIA | R | 2.11 |
| 4 | apache1 | DIA | - | DIA | DIA | R | 1.94 |
| 5 | bfly | ELL | - | ELL | ELL | R | 1.18 |
| 6 | whitaker3_dual | ELL | - | ELL | ELL | R | 4.89 |
| 7 | ch7-9-b3 | ELL | - | ELL | ELL | R | 2.25 |
| 8 | shar_te2-b2 | ELL | - | ELL | ELL | R | 2.24 |
| 9 | pkustk14 | $confidence < TH$ | CSR+COO | CSR | CSR | W | 16.39 |
| 10 | crankseg_2 | $confidence < TH$ | CSR+COO | CSR | CSR | W | 16.28 |
| 11 | Ga3As3H12 | $confidence < TH$ | CSR+COO | CSR | CSR | W | 16.2 |
| 12 | HV15R | $confidence < TH$ | CSR+COO | CSR | CSR | W | 15.43 |
| 13 | europe_osm | COO | - | COO | COO | R | 2.3 |
| 14 | D6-6 | COO | - | COO | COO | R | 5.79 |
| 15 | dictionary28 | COO | - | COO | COO | R | 2.05 |
| 16 | roadNet-CA | COO | - | COO | COO | R | 2.38 |

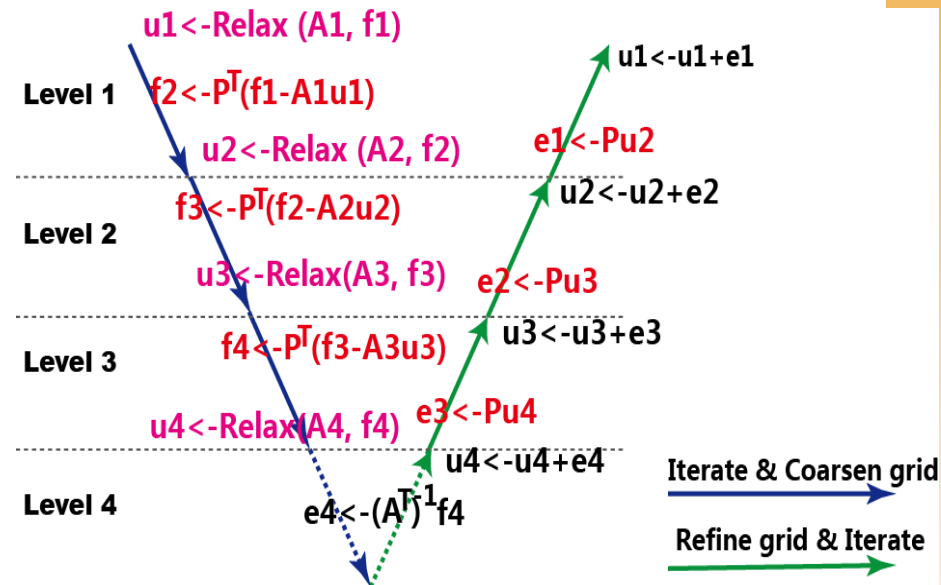"R" and "W" represent Right and Wrong prediction respectively.

## ◆ Algebraic Multi-grid Algorithm

- An iterative algo. to solve linear equations $Au=f$, where A is sparse matrix, u, f are dense vectors
- As a pre-conditioner applied in applications such as laser fusion and climate modeling

## ◆ SpMV the critical operation of AMG, takes 90% execution time.

| Coarsen | Rows | Hypre AMG | SMAT AMG | Speedup |
|---|---|---|---|---|
| cljp_7pt_50 | 125k | 3034 | 2487 | **1.22** |
| rugeL_9pt_500 | 250k | 388 | 300 | **1.29** |

**AMG Execution Process**



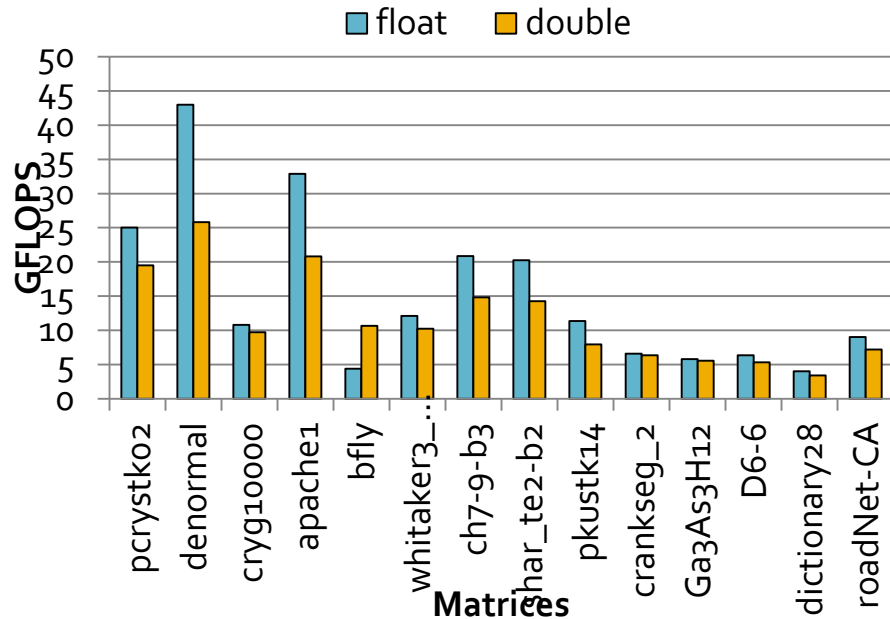| Level 1 | $u1 \leftarrow Relax\ (A1, f1)$ <br> $f2 \leftarrow -P^T(f1-A1u1)$ | $u1 \leftarrow u1+e1$ |
| Level 2 | $u2 \leftarrow Relax\ (A2, f2)$ <br> $f3 \leftarrow -P^T(f2-A2u2)$ | $e1 \leftarrow Pu2$ <br> $u2 \leftarrow u2+e2$ |
| Level 3 | $u3 \leftarrow Relax(A3, f3)$ <br> $f4 \leftarrow -P^T(f3-A3u3)$ | $e2 \leftarrow Pu3$ <br> $u3 \leftarrow u3+e3$ |
| Level 4 | $u4 \leftarrow Relax(A4, f4)$ <br> $e4 \leftarrow -(A^T)^{-1}f4$ | $e3 \leftarrow Pu4$ <br> $u4 \leftarrow u4+e4$ |

Iterate & Coarsen grid

Refine grid & Iterate
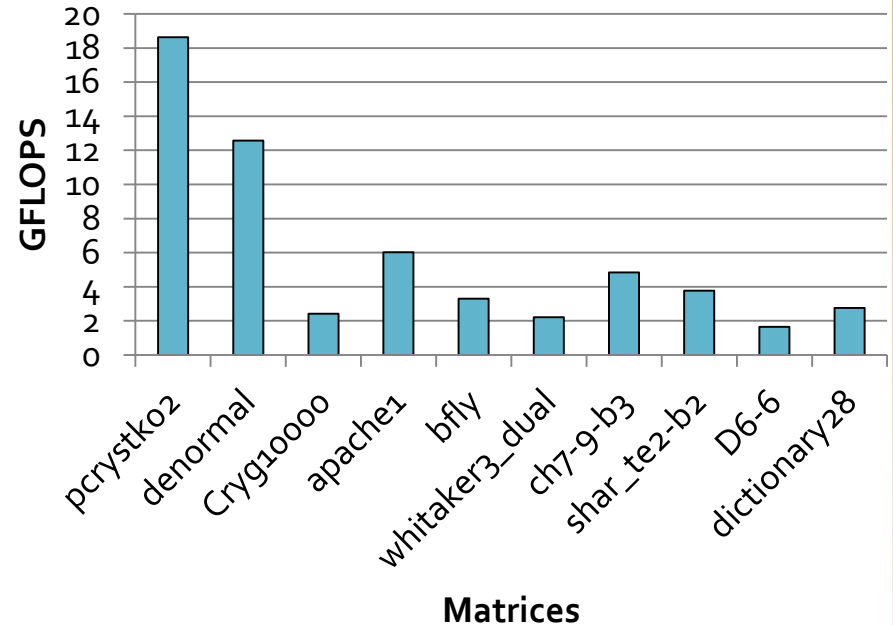
- Relax is relaxation algotithm, such as Jacobi、G-S iterative methods
- A, P are sparse matrices
- U, f, e are dense vectors

# SMAT Many-core

## Performance



NVIDIA K20



Xeon Phi

**Accuracy**
- For 289 testing matrices: 89%(single), 95%(double)

# SMAT Summary

◆ **Develop auto-tuning method**

- Design application-architecture aware SpMV auto-tuner.
- Develop auto-tuning method to algorithm level

◆ **Introduce data mining to auto-tuning method**

- Reinforce its usability and extensibility

◆ **API Easy-to-use**

- Unify the interface

◆ **Increase SpMV performance using SMAT**

◆ **Apply SMAT to AMG, and extend it to many-core architecture**

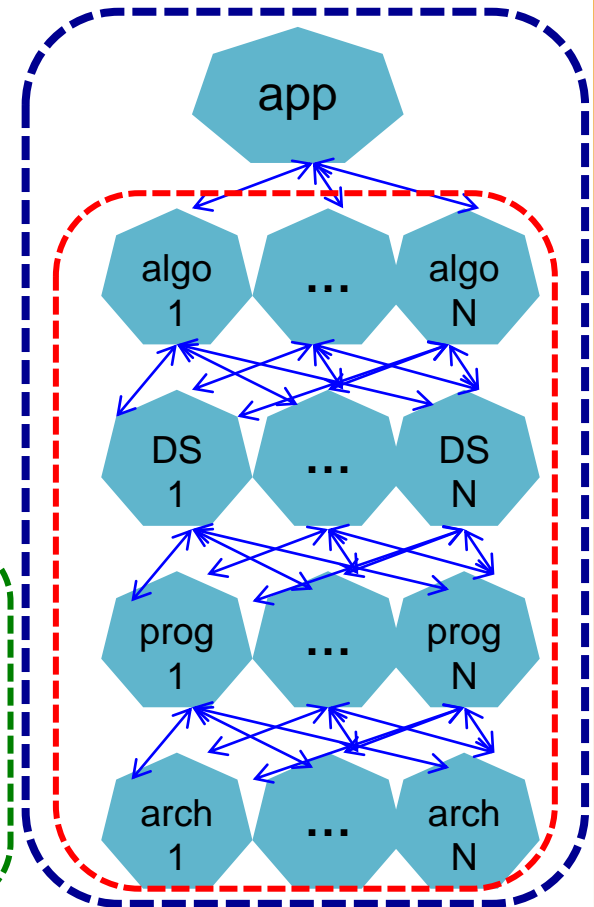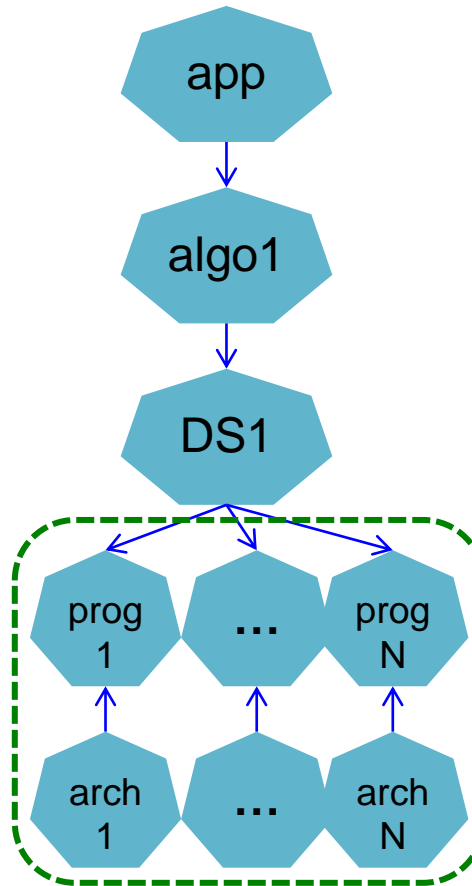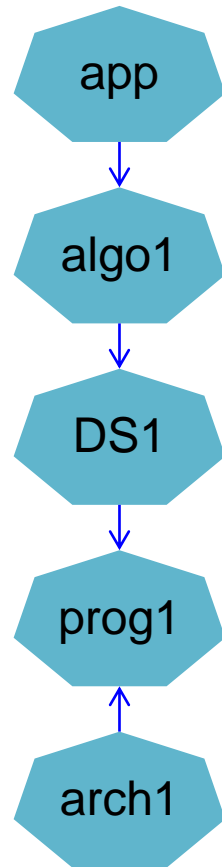# High Performance Computational Software Development

◆**Extend storage formats**

◆**Combine other auto-tuners**

Offline → Online ⟶ ······▶ Call

Application

Architecture

Extract App. Feature

Extract Arch. Feature

Parameter Set

Arch.-aware Auto-tuner (e.g. OSKI)

Get Para. Value (Learning Set)

1st level Auto-tuning

Feature Database

Best Implementation

Summarize Features

Model

Predict

Best Algo. & Impl.

Extract Features

Application

Online

# Thank You !

Question?

SMAT: An Input Adaptive Auto-Tuner
for Sparse Matrix-Vector Multiplication.
Jiajia Li, Guangming Tan, Mingyu Chen, Ninghui Sun

**ICT 中国科学院计算技术研究所**
Institute of Computing Technology, Chinese Academy of Sciences