



# SPOOF: Sum-Product Optimization and Operator Fusion for Large-Scale Machine Learning

**Tarek Elgamal<sup>2</sup>, Shangyu Luo<sup>3</sup>, Matthias Boehm<sup>1</sup>,  
Alexandre V. Evfimievski<sup>1</sup>, Shirish Tatikonda<sup>4</sup>,  
Berthold Reinwald<sup>1</sup>, Prithviraj Sen<sup>1</sup>**

<sup>1</sup> IBM Research – Almaden

<sup>2</sup> University of Illinois

<sup>3</sup> Rice University

<sup>4</sup> Target Corporation

# Motivation

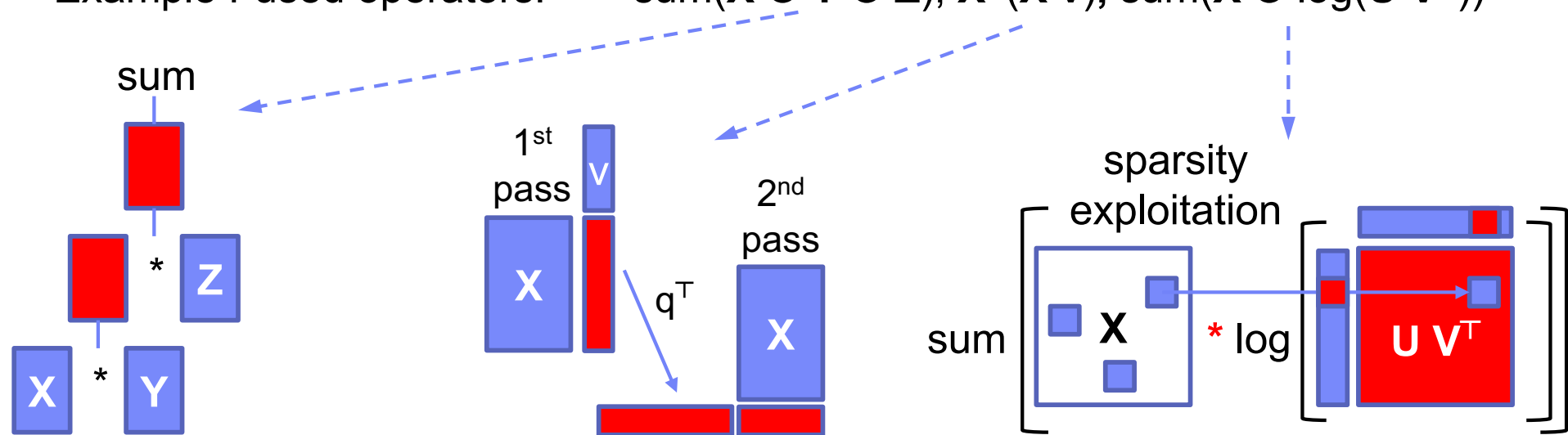
## ■ Declarative Large-Scale Machine Learning (ML)

- Simplify development / usage of ML tasks or algorithms
- SystemML: High-level language  $\rightarrow$  data independence / plan generation
- State-of-the-art compilers: **rewrites, operator selection, fused operators**



## ■ Ubiquitous Optimization Opportunities

- Example Rewrites:  $\mathbf{X}^T \mathbf{y} \rightarrow (\mathbf{y}^T \mathbf{X})^T$ ,  $\text{sum}(\lambda \mathbf{X}) \rightarrow \lambda \text{sum}(\mathbf{X})$ ,  
 $\text{trace}(\mathbf{X} \mathbf{Y}) \rightarrow \text{sum}(\mathbf{X} \odot \mathbf{Y}^T)$
- Example Fused operators:  $\text{sum}(\mathbf{X} \odot \mathbf{Y} \odot \mathbf{Z})$ ,  $\mathbf{X}^T (\mathbf{X} \mathbf{v})$ ,  $\text{sum}(\mathbf{X} \odot \log(\mathbf{U} \mathbf{V}^T))$



# Motivation

## ■ Declarative Large-Scale Machine Learning (ML)

- Simplify development / usage of ML tasks or algorithms
- SystemML: High-level language  $\rightarrow$  data independence / plan generation
- State-of-the-art compilers: **rewrites, operator selection, fused operators**



## ■ Ubiquitous Optimization Opportunities

- Example Rewrites:  $\mathbf{X}^T \mathbf{y} \rightarrow (\mathbf{y}^T \mathbf{X})^T$ ,  $\text{sum}(\lambda \mathbf{X}) \rightarrow \lambda \text{sum}(\mathbf{X})$ ,  
 $\text{trace}(\mathbf{X} \mathbf{Y}) \rightarrow \text{sum}(\mathbf{X} \odot \mathbf{Y}^T)$
  - Example Fused operators:  $\text{sum}(\mathbf{X} \odot \mathbf{Y} \odot \mathbf{Z})$ ,  $\mathbf{X}^T (\mathbf{X} \mathbf{v})$ ,  $\text{sum}(\mathbf{X} \odot \log(\mathbf{U} \mathbf{V}^T))$
- $\rightarrow$  Fewer intermediates, fewer scans, sparsity exploitation, less compute**

## ■ Problems and Challenges

- **Large Development Effort:** number of patterns, multiple runtime back-ends, multiple formats and combinations (sparse/dense)
- **High Performance Impact:** slightly changed patterns can render rewrites and fused operators inapplicable

# Example PNMF – A 1000x War Story

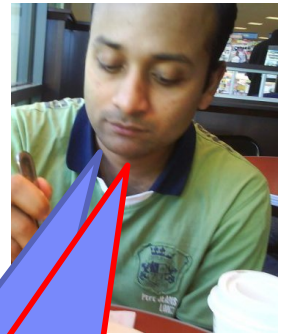
## ■ Poisson Nonnegative Matrix Factorization (PNMF)

–  $X \approx W H$  of low rank  $k$ ;  $X$ : 200K x 200K, sp=0.001 (480MB)

```

1: X = read("./input/X")
2: k = 100; eps = 1e-15; max_iter = 10; iter = 1;
3: W = rand(rows=nrow(X), cols=k, min=0, max=0.025)
4: H = rand(rows=k, cols=ncol(X), min=0, max=0.025)
5: while( iter < max_iter ) {
6:   H = (H*(t(W)%*(X/(W%*H+eps)))) / t(colSums(W));
7:   W = (W*((X/(W%*H+eps))%*t(H))) / t(rowSums(H));
8:   obj = sum(W%*H) - sum(X*log(W%*H+eps));
9:   print("iter=" + iter + " obj=" + obj);
10:  iter = iter + 1;
11: }
12: write(W, "./output/W");
13: write(H, "./output/H");

```



It still takes forever  
– btw, I changed it slightly

Here is an interesting  
rewrite:  $\text{sum}(W H) \rightarrow$   
 $\text{colSums}(W) \text{ rowSums}(H)$

The problem is  $W H$  (320GB), but  
we could add sparsity-exploiting  
**fused operators and rewrites**

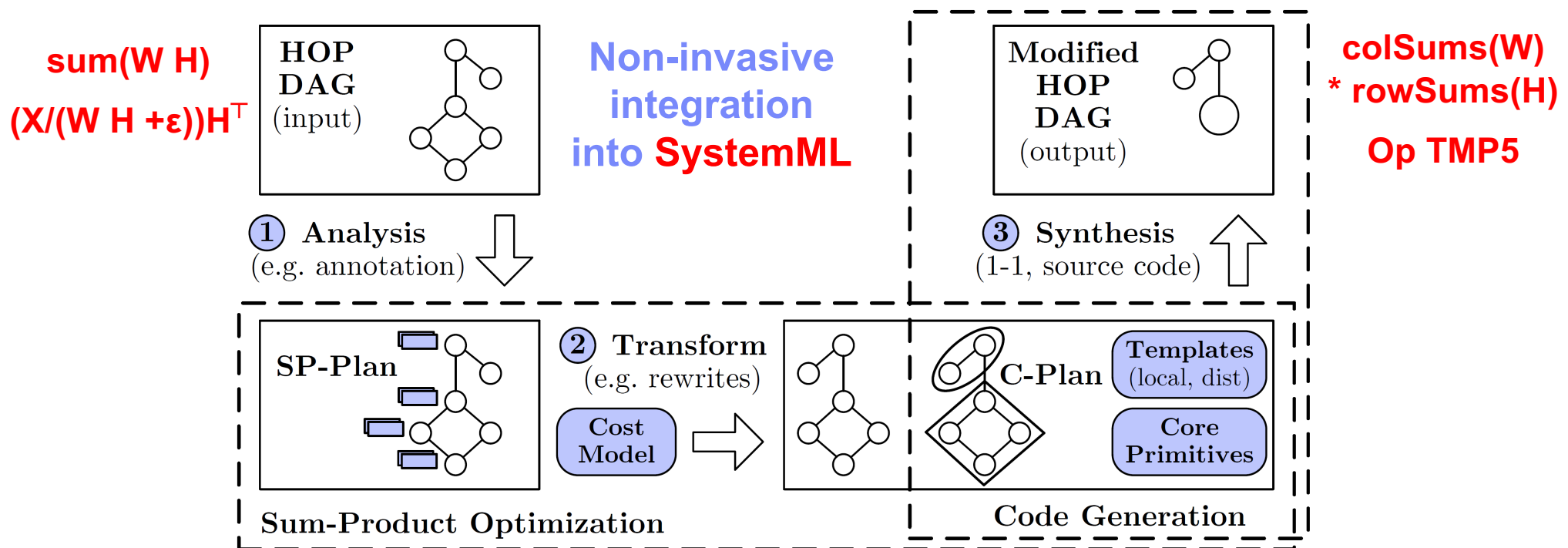
➔ rewrites and fused operators: **1000x**



# Our Vision: Holistic Optimization Framework

## ■ SPOOF Compiler Framework

- Automatic rewrite identification and operator fusion
- Increased opportunities and side effects (CSE, rewrites  $\leftrightarrow$  fusion)
- **Key ideas:** (1) break up LA operations into basic operators (in RA),  
(2) elementary sum-product and RA rewrites, and (3) fused operator generation





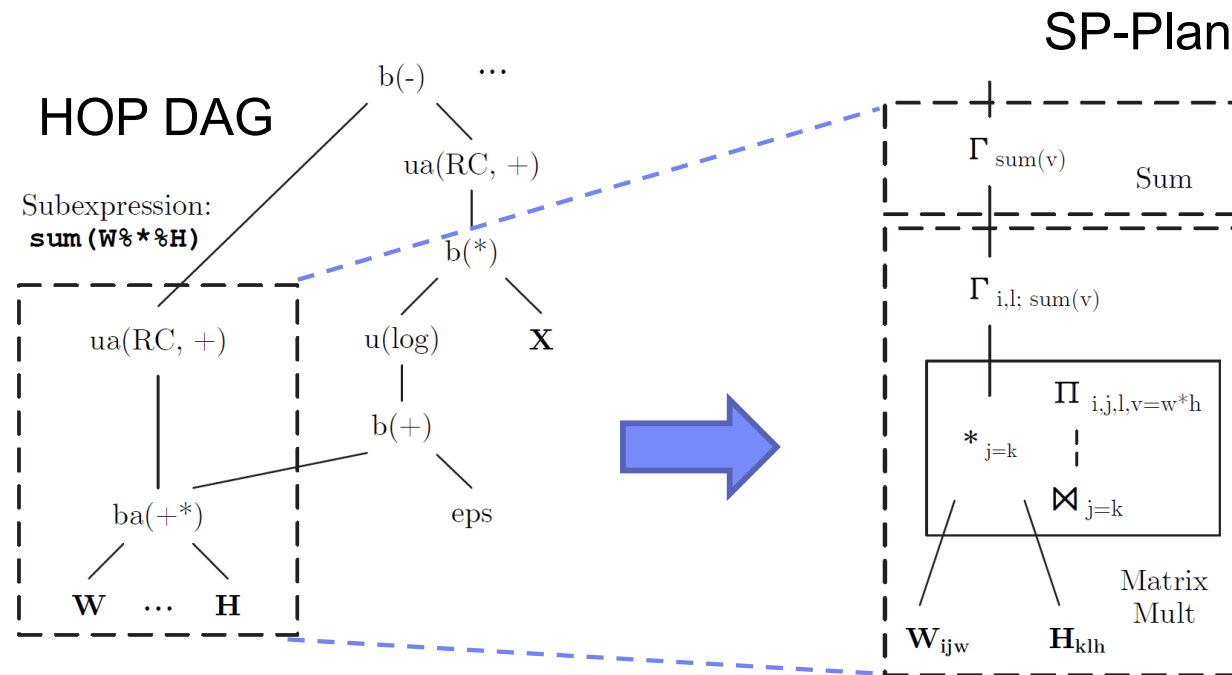
# Sum-Product Optimization

## ■ SP-Plan Representation: restricted relational algebra

- **Data:** input matrices are relations of  $(i, j, v)$ -tuples (intermediates are tensors)
- **Basic operations:** selection  $\sigma$ , extended projection  $\Pi$ , aggregation  $\Gamma$ , join  $\bowtie$
- **Composite operations:** e.g., multiply  $A_{ij} *_{i=k \wedge j=l} B_{kl} := \Pi_{i,j;a*b}(A_{ija} \bowtie_{i=k \wedge j=l} B_{klb})$   
addition  $A_{ij} +_{i=k \wedge j=l} B_{kl} := \Pi_{i,j;a+b}(A_{ija} \bowtie_{i=k \wedge j=l} B_{klb})$
- **Two restrictions:** a single value attribute per relation, and unique composite indexes per relation → **single value per tensor cell**

## ■ Example SP Plan

- **sum(W H)**

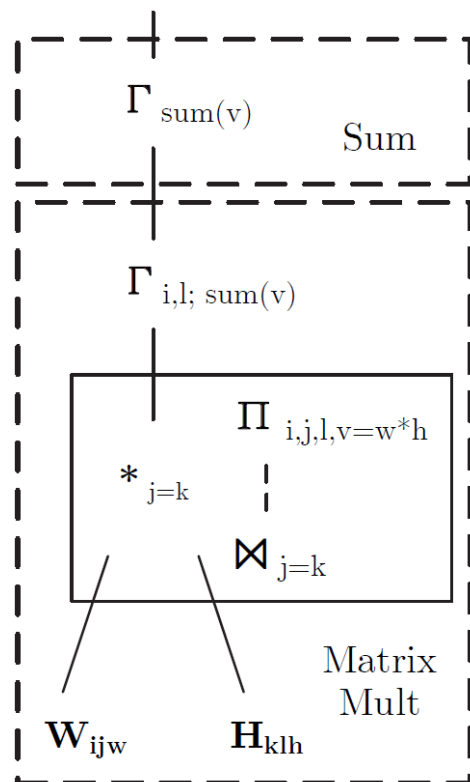


# Sum-Product Optimization, cont.

## Example SP Plan Rewrites

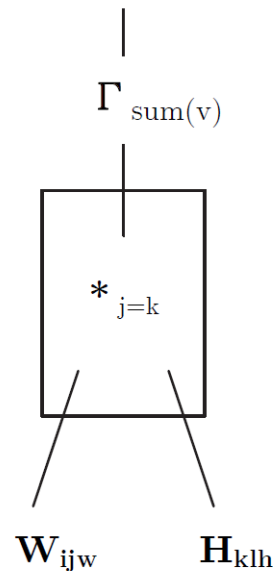
—  $W := 200K \times 100$ ,  $H := 100 \times 200K$

**But, SP opt alone can be counter-productive (e.g., CSE 'W H')**



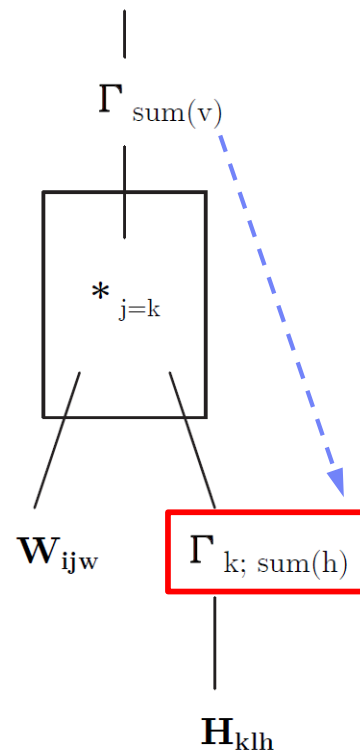
**8.04 TFLOPs**

### Aggregation Elimination



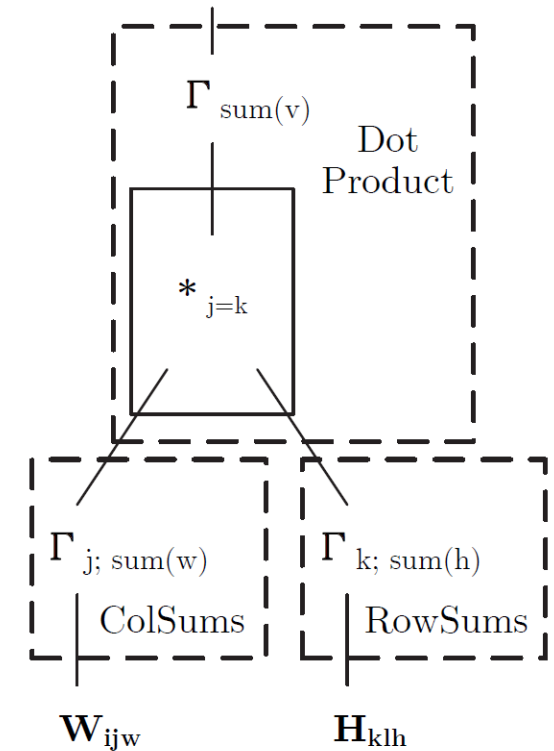
**8 TFLOPs**

### Sum-Product (distributive law)



**60 MFLOPs**

### Sum-Product (distributive law)



**40 MFLOPs**

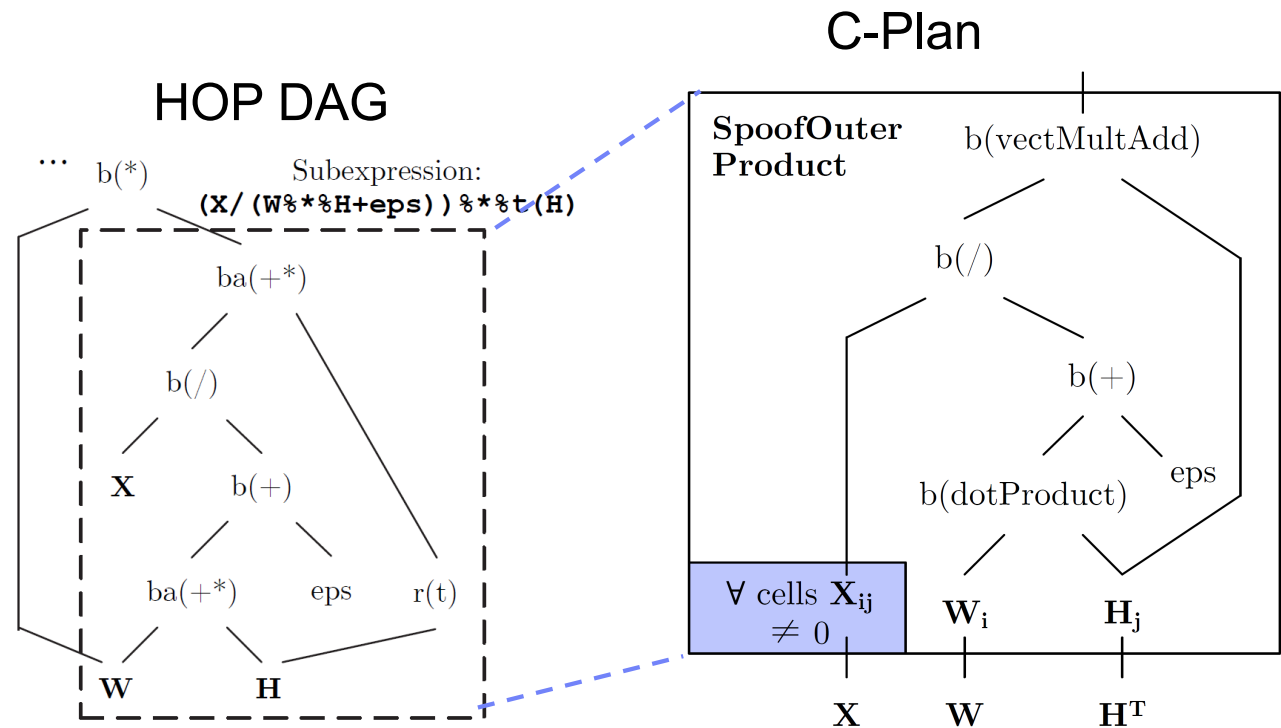
# Operator Fusion

## ■ C-Plan Representation

- **Hybrid approach:** hand-coded operator skeletons with custom body code  
 ➔ **Efficiency (data access, multi-threading) and flexibility**
- **Template C-Nodes:** generic fused operator skeletons (w/ data binding)  
 e.g., SpoofOuterProduct, SpoofCellwise, SpoofRowAggregate
- **Primitive C-Nodes:** vector/scalar operations

## ■ Example C-Plan

- $(X / (W H + \epsilon)) H^T$   
 (PNMF update rule)

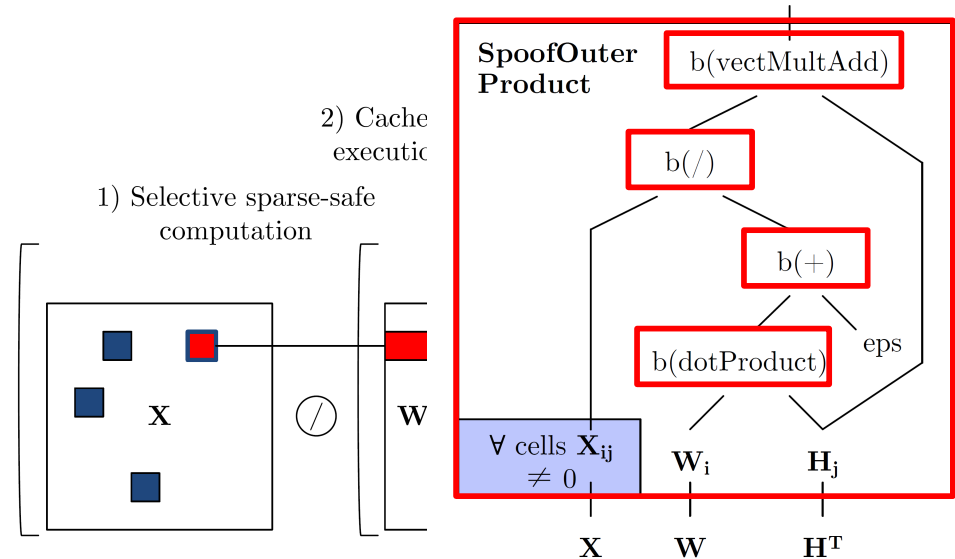




# Operator Fusion, cont.

## ■ Example C-Plan Codegen

- Recursive codegen on C-Plan
- **Generated operator inherits data access, multi-threading, etc from template skeleton**



```

1: public final class TMP5 extends SpoofOuterProduct {
2:     public TMP5() {
3:         _type = OuterProductType.RIGHT;
4:     }
5:     protected void exec(double a, double[] b, int bi,
6:         double[] c, int ci, ..., double[] d, int di, int k)
7:     {
8:         double TMP1 = dotProduct(b, c, bi, ci, k);           // WH
9:         double TMP2 = TMP1 + 1.0E-15;                       // +eps
10:        double TMP3 = a / TMP2;                               // X/
11:        vectMultiplyAdd(TMP3, c, d, ci, di, k);               // t(H)
12:    }
13: }

```

**Custom  
body  
code**

# Experimental Setting

## ■ Cluster Setup

- 1 head node (2x4 Intel E5530, 64GB RAM), and  
6 worker nodes (2x6 Intel E5-2440, 96GB RAM, 12x2TB disks)
- Spark 1.5.2 with 6 executors (24 cores, 60GB), 30GB driver memory

## ■ ML Programs and Data

- 3 full-fledged ML algorithms (PNMF, L2SVM, Mlogreg)
- Synthetically generated data

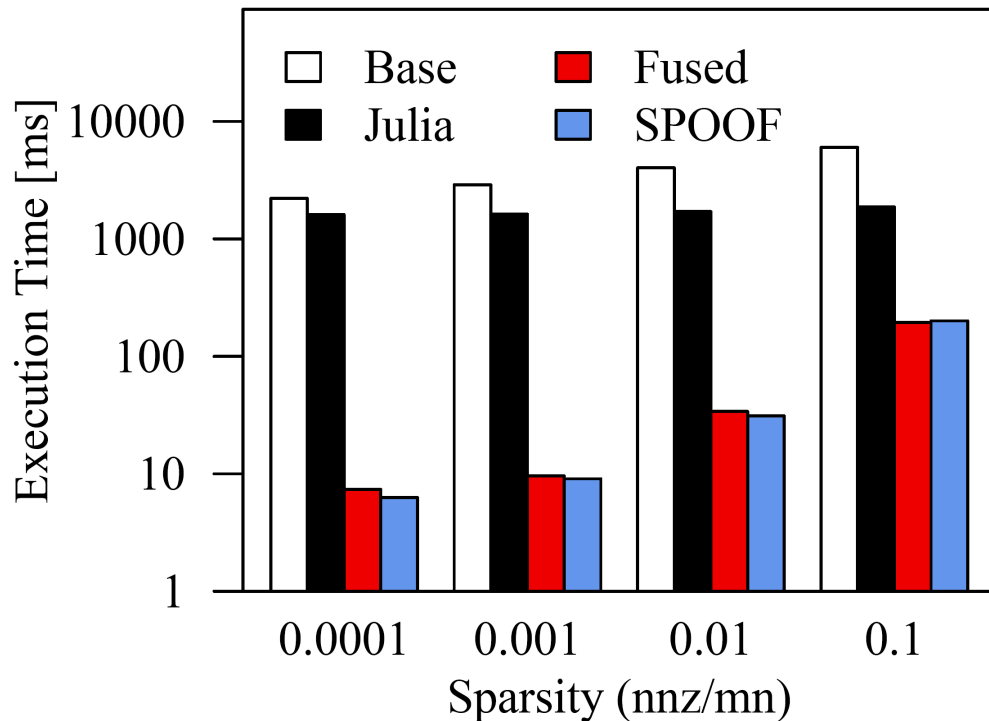
## ■ Selected Baselines

- Apache SystemML 0.10 (May 2016): **Base, Fused, SPOOF**
- **Julia** 0.5 (Sep 2016) w/ LLVM-based just-in-time compiler

# Micro Benchmarks: Operations Performance

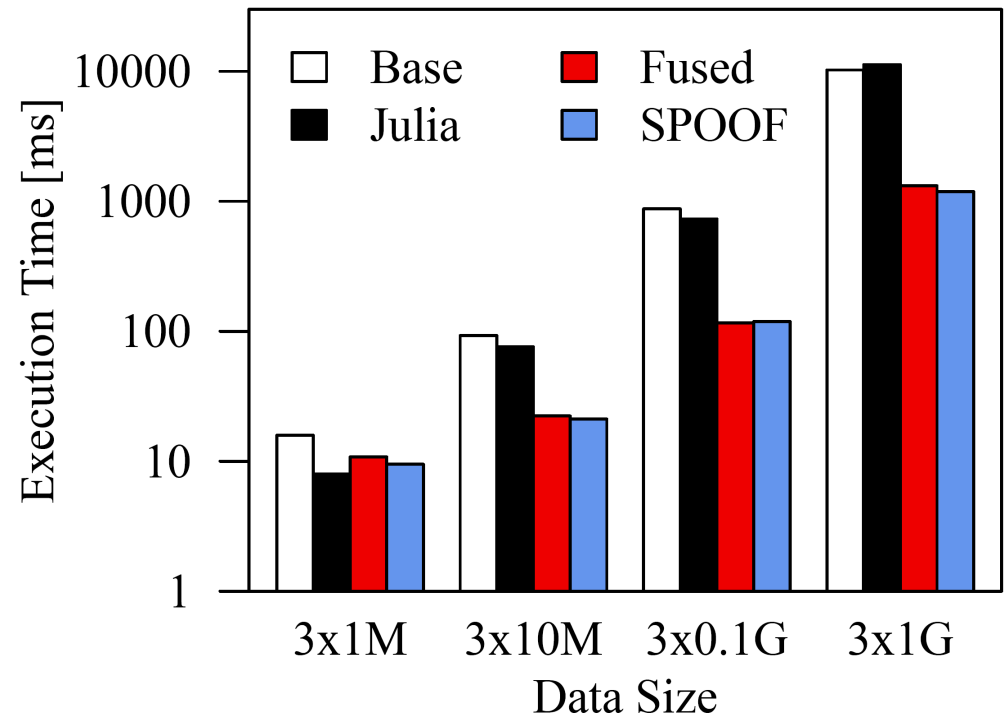
(@ single worker node)

$(\mathbf{X} / (\mathbf{W} \mathbf{H} + \epsilon)) \mathbf{H}^T$ , (PNMF)  
10K x 10K, k=100, Multi-threaded



→ Sparsity-exploiting operators  
at 1/12 peak compute bandwidth

$\text{sum}(\mathbf{X} \odot \mathbf{Y} \odot \mathbf{Z})$ , (L2SVM)  
dense, Multi-threaded



→ Fused operator w/o  
intermediates at peak 1xlocal /  
remote memory bandwidth (25GB/s)

# End-to-End Experiments: PNMf and LSVM

## ■ PNMf Execution Time (incl. compilation and I/O)

- 20 iterations, rank  $k = 100$

Dataset	Base	Fused	SPOOF
10K x 10K, 0.001	<b>251 s</b>	<b>6 s</b>	9 s
25K x 25 K, 0.001	<b>4,748 s</b>	<b>9 s</b>	11 s
200K x 200K, 0.001	<b>&gt;24h</b>	<b>121 s</b>	125 s

## ■ L2SVM Execution Time (incl. compilation and I/O)

- 20 outer iterations,  $\epsilon = 10^{-14}$

Dataset	Base	Fused	SPOOF
100K x 10, 1.0 (8MB)	3 s	<b>3 s</b>	5 s
1M x 10, 1.0 (80MB)	9 s	<b>7 s</b>	8 s
10M x 10, 1.0 (800MB)	50 s	34 s	<b>17 s</b>
100M x 10, 1.0 (8GB)	525 s	320 s	<b>114 s</b>

# Conclusions

## ■ Summary

- **SPOOF: Automatic rewrite identification and operator fusion**
- Non-invasive compiler/runtime integration into SystemML
- Plan representation/compilation for sum-product and codegen

## ■ Conclusions and Future Work

- **Many rewrite/fusion opportunities with huge performance impact**
- Performance close to hand-coded ops w/ moderate compilation overhead
- Future work: distributed operations, optimization algorithms

## ■ Available Open Source (soon)

- SYSTEMML-448: Code Generation, experimental in 1.0 release
- Sum-product optimization and fusion optimizations later





**SystemML is Open Source:**

Apache Incubator Project since 11/2015

Website: <http://systemml.apache.org/>

Sources: <https://github.com/apache/incubator-systemml>

# Backup: Operator Fusion, cont.

## ■ Example C-Plan Codegen

### – L2SVM inner loop

```

1: public final class TMP2 extends SpoofCellwise {
2:   public TMP2() {
3:     _type = CellType.FULL_AGG;
4:   }
5:   protected double exec(double a, double[][]
6:     vectors, double[] scalars, ..., int rowIndex)
7:   {
8:     double TMP3 = vectors[1][rowIndex];
9:     double TMP4 = vectors[0][rowIndex];
10:    double TMP5 = a * scalars[0];
11:    double TMP6 = TMP4 + TMP5;
12:    double TMP7 = TMP3 * TMP6;
13:    double TMP8 = 1 - TMP7;
14:    double TMP9 = (TMP8 > 0) ? 1 : 0;
15:    double TMP10 = TMP8 * TMP9;
16:    double TMP11 = TMP10 * TMP3;
17:    double TMP12 = TMP11 * a;
18:    return TMP12;
19:  }
20: }

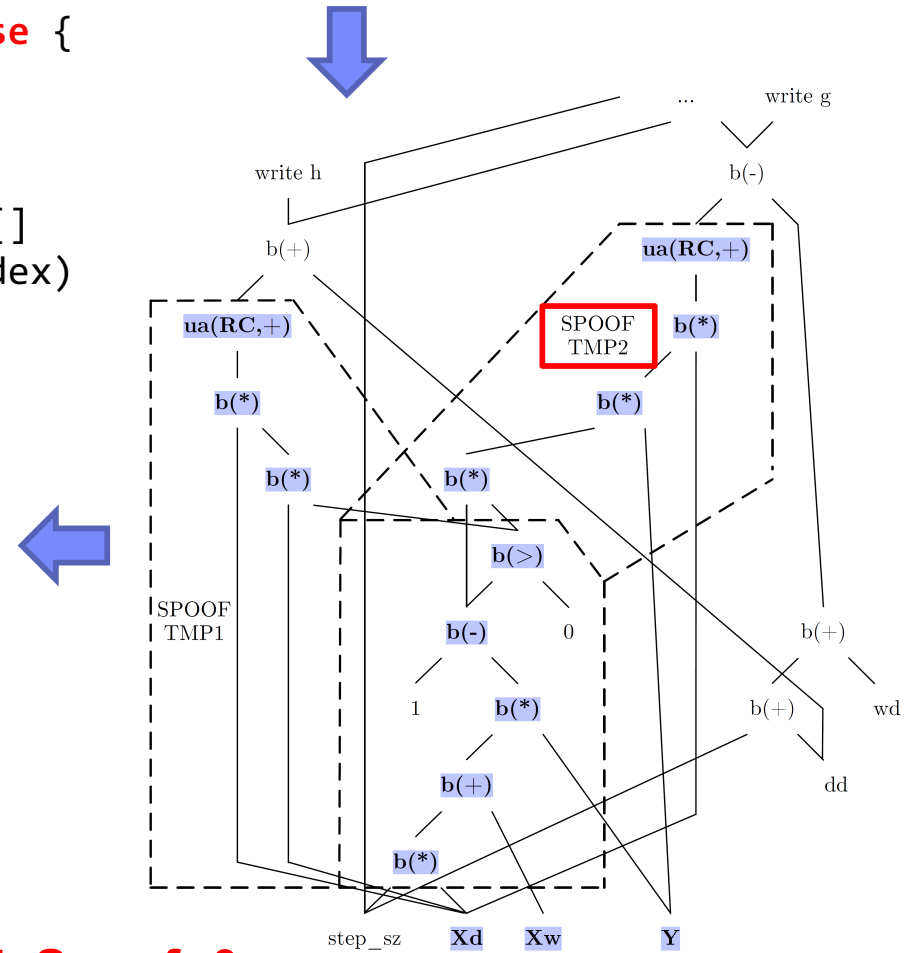
```

# Intermediates: Base: 10, Fused: 5, Spoof: 0

```

1: out = 1 - Y * (Xw + step_sz*Xd);
2: sv = (out > 0);
3: out = out * sv;
4: g = wd+step_sz*dd - sum(out*Y*Xd);
5: h = dd + sum(Xd*sv*Xd);
6: step_sz = step_sz - g/h;

```



# Backup: Plan Caching Effects for Mlogreg

## ■ Dynamic Recompilation

- Problem of unknown or changing sizes (e.g., UDFs, data-dep. ops, size expr.)
- Integration of Spoof into dynamic recompiler → **huge compilation overhead**
- **Plan cache**: reuse compiled ops across DAGs / recompilations

## ■ Mlogreg Cache Statistics

- 500K x 200 (800MB), 20/5 outer/inner iterations,  $\epsilon = 10^{-14}$
- CSLH: Context-sensitive literal heuristic

Statistic	Spoof no PC	Spoof constant PC	Spoof CSLH
Execution time	<b>49.29 s</b>	19.87 s	<b>14.48 s</b>
PC hit rates	0 / 462	388 / 462	<b>449 / 462</b>
Javac compile time (sync)	<b>34.45 s</b>	6.88 s	<b>1.97 s</b>
JIT compile time (async)	25.36 s	18.84 s	<b>10.50 s</b>