

MNC: Structure-Exploiting Sparsity Estimation for Matrix Expressions

Johanna Sommer³, Matthias Boehm², Alexandre Evfimievski¹, Berthold Reinwald¹, Peter Haas⁴

SIGMOD 2019
Research 16: Machine Learning

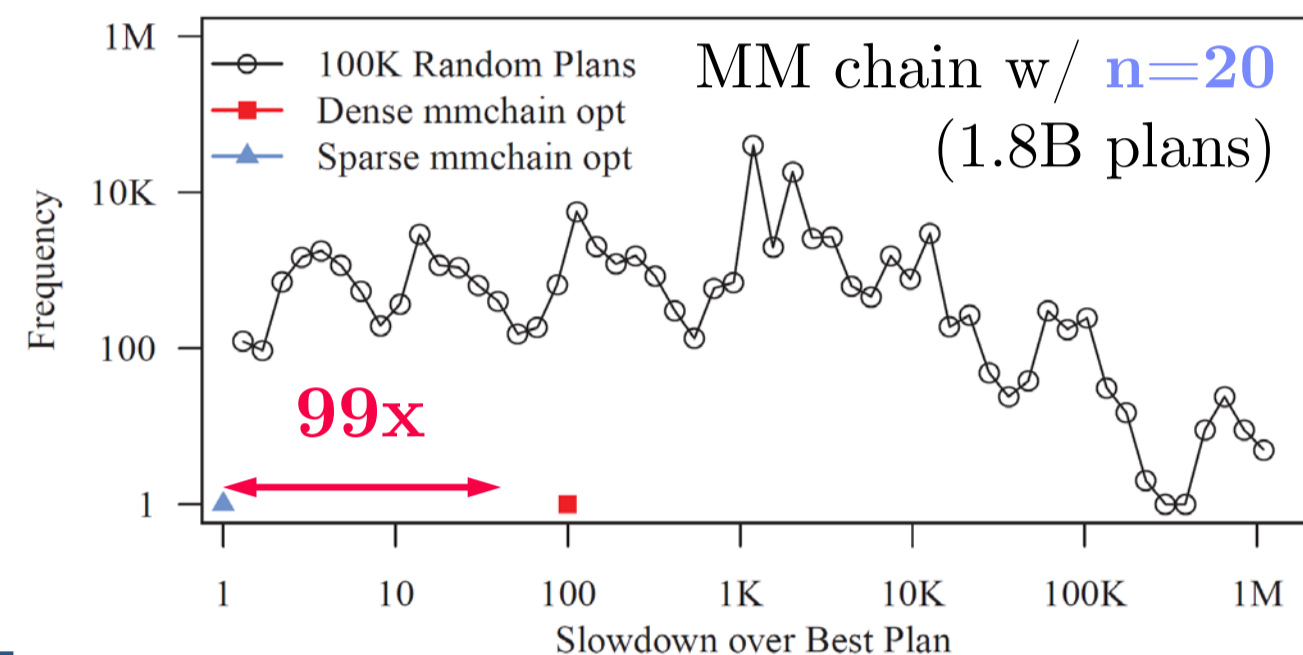
¹ IBM Research – Almaden ² Graz University of Technology
³ IBM Germany ⁴ UMass Amherst

Contact: Matthias Boehm
m.boehm@tugraz.at

Motivation

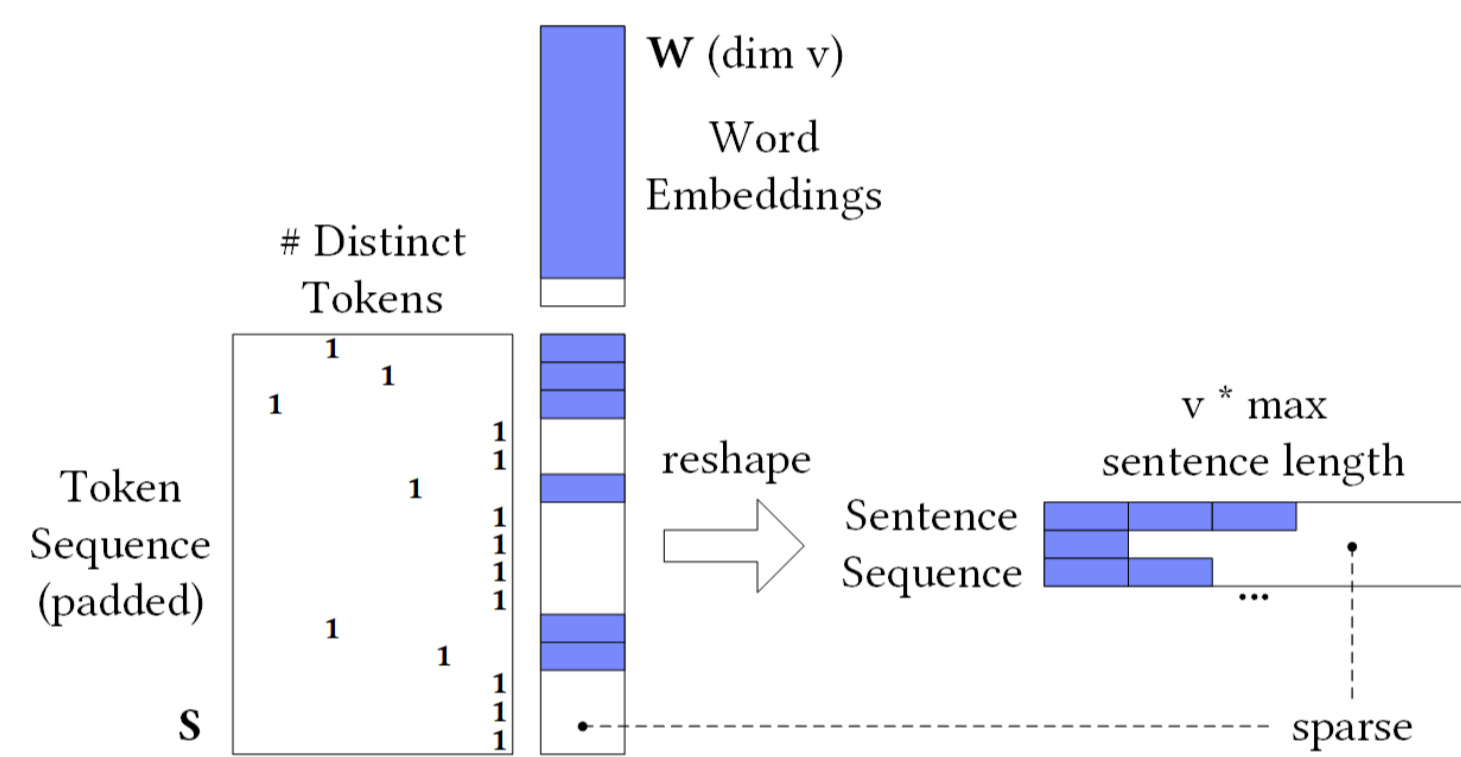
Sparsity Estimation

- Ubiquitous Sparse Matrices
 - Sparse inputs: NLP, graph, recommenders, scientific computing
 - Data preprocessing: binning/recoding + one-hot encoding/feature hashing
 - Sparse intermediates: selection predicates, dropout
 - Transformation matrices: random reshuffling, sampling
 - Automatic application of sparse operations
- Problem of Sparsity Estimation
 - Compilation: memory and cost estimation (local vs distributed, rewrites, resource allocation, operator fusion)
 - Runtime: format decisions and memory pre-allocation



Example: NLP Sentence Encoding

- Problem
 - Encoding a sequence of words (padded to max sentence length) into sequence of word embeddings
 - Input for SentenceCNN
- Implementation
 - Permutation matrix multiply
 - Reshape into sentences



Sparsity Estimators

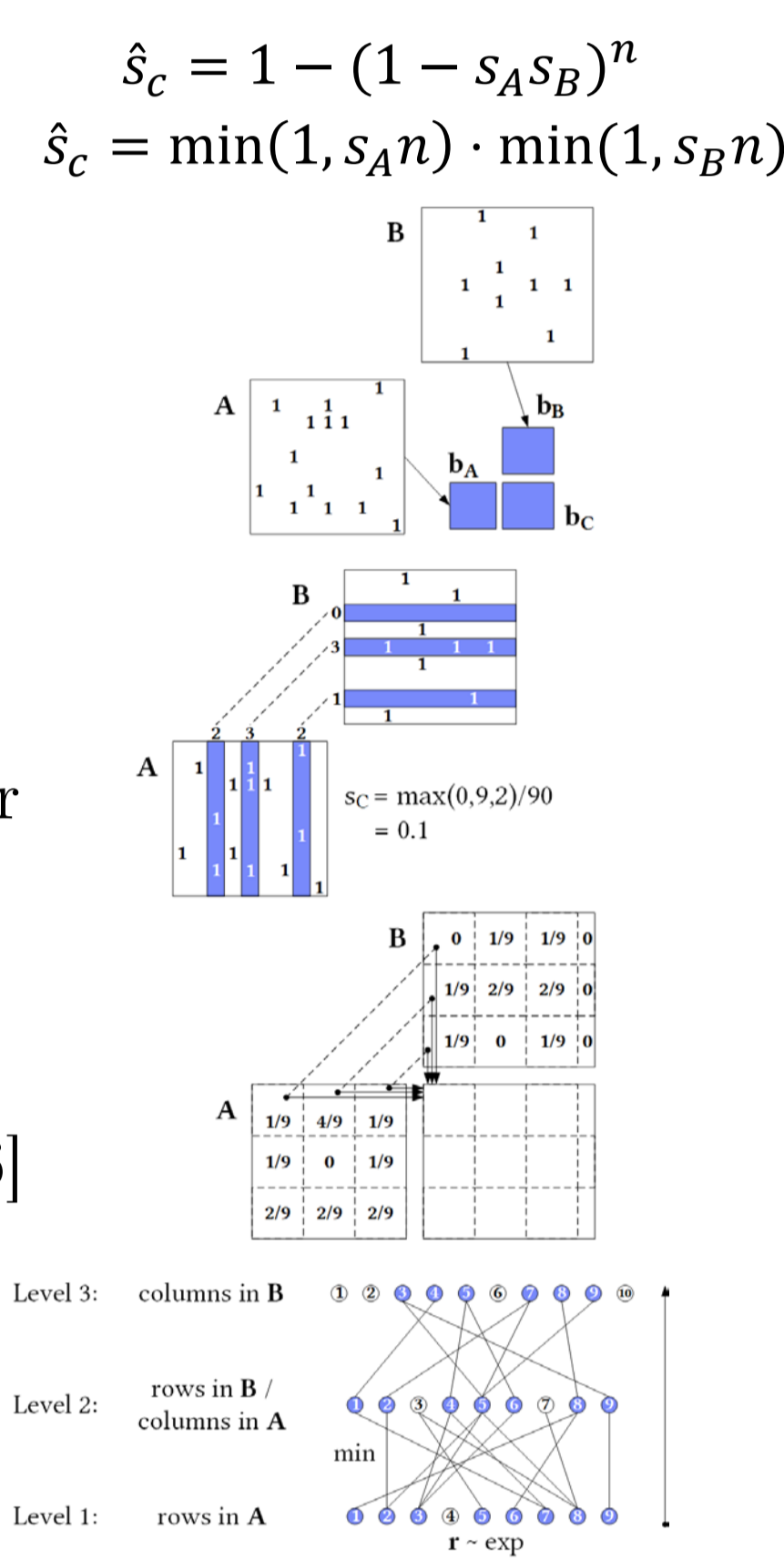
Assumptions

- A1: No Cancellation Errors
- A2: No Not-A-Number (NaN*0=NaN)

Common assumptions
→ Boolean mat. mult.
 $m \times n$ by $n \times l$
 $(s_C = \frac{nnz(C)}{ml})$

Existing Sparsity Estimators

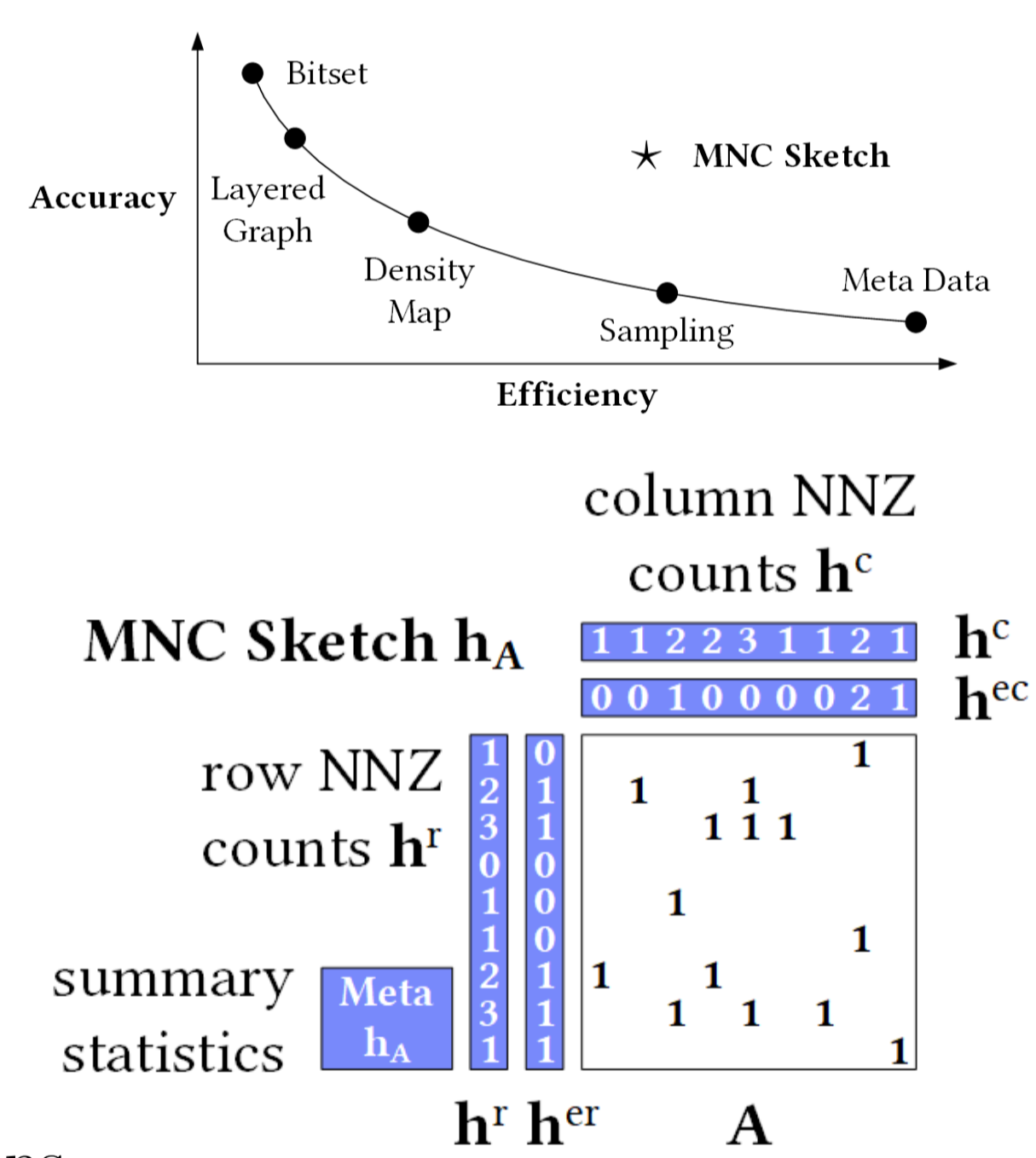
- Naïve Metadata Estimators**
 - Estimate sparsity from the inputs (e.g., SystemML)
- Naïve Bitset Estimator**
 - Convert inputs to bitsets and perform Boolean mm
 - Examples: SciDB [SSDBM'11], cuSparse, MKL
- Sampling**
 - Sampling of aligned columns of A and rows of B
 - Sparsity estimated via max of count-products
 - Examples: MatFast [ICDE'17], improvements in paper
- Density Map**
 - Store sparsity per $b \times b$ block (default $b = 256$)
 - MM-like estimator (E_{ac} for $*$, $s_A + s_B - s_A s_B$ for $+$)
 - Examples: SpMacho [EDBT'15], AT Matrix [ICDE'16]
- Layered Graph [J.Comb.Opt.'98]**
 - Nodes: rows/columns in mm chain
 - Edges: non-zeros connecting rows/columns
 - Assign r-vectors $\sim \exp$ and propagate via min



MNC Sketch

Basic Data Structure

- Row/Column NNZs**
 - Count vectors for nnz per row/col
 - $h^r = \text{rowSums}(A \neq 0)$; $h^c = \text{colSums}(A \neq 0)$
- Extended Row/Column NNZs**
 - $h^{er} = \text{rowSums}((A \neq 0) * (h^c = 1))$
 - $h^{ec} = \text{colSums}((A \neq 0) * (h^r = 1))$
- Summary Statistics**
 - Max nnz per row/column
 - # of non-empty / # of half-full rows/columns



+ Additional Operations
(reorg and element-wise ops in the paper)

Sparsity Estimation

- Basic Estimator
 - DensityMap-like estimator over column/rows slices
- Exact Sparsity Estimation
 - Under assumptions A1 and A2, exact sparsity estimation possible

$$s_C = \hat{s}_C = \frac{h_A^c h_B^r}{(ml)}$$

if $\max(h_A^r) \leq 1 \vee \max(h_B^c) \leq 1$

Sparsity Estimation, cont.

- Lower and Upper Bounds**
 - Lower bound → clipping of bad estimates
 - Upper bound → improved estimates
- Extended NNZ Counts**
 - Compose estimate from exactly known and estimated quantities (i.e., number non-zeros)

$$s_C \geq |h_A^r > n/2| \cdot |h_B^c > n/2| / (ml)$$

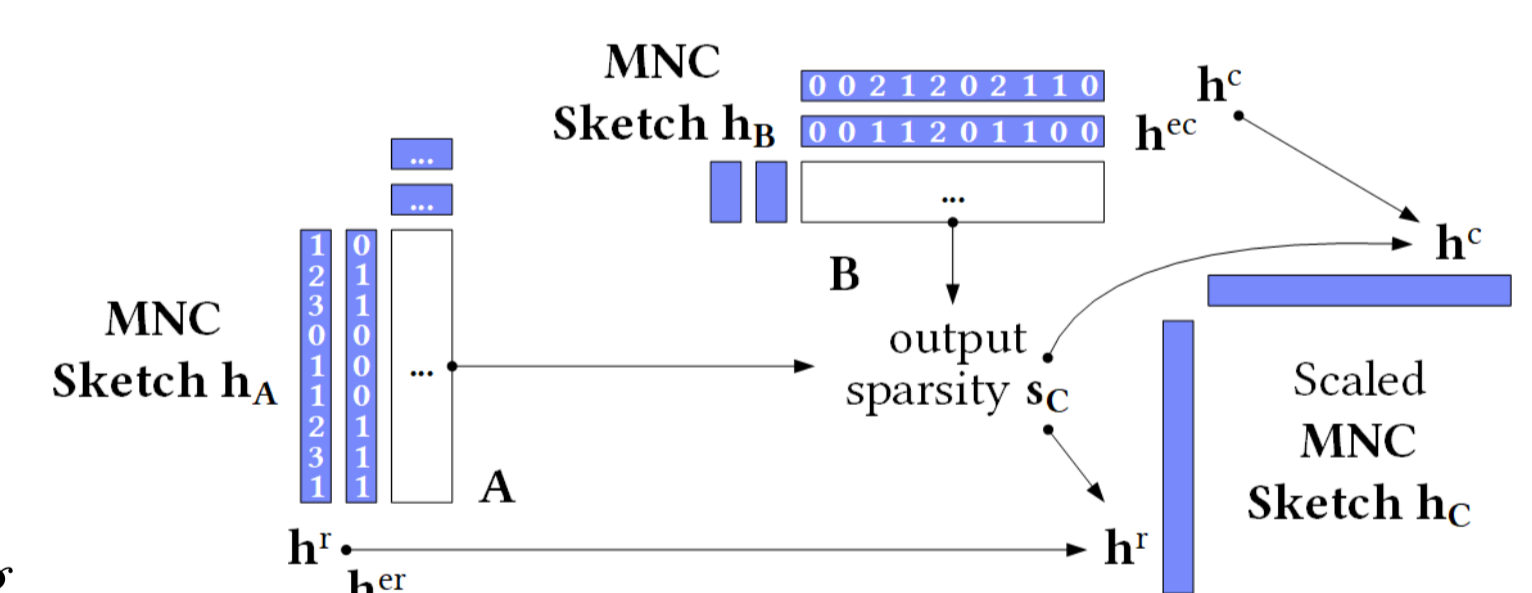
$$s_C < nnz(h_A^r) \cdot nnz(h_B^c) / (ml)$$

$$\hat{s}_C = (h_A^{ec} h_B^r + (h_A^c - h_A^{ec}) h_B^{er}) / (ml) + E_{dm}(h_A^c - h_A^{ec}, h_B^r - h_B^{er}, p) \cdot p / (ml)$$

$$p = \frac{(nnz(h_A^r) - |h_A^r = 1|)}{(nnz(h_B^c) - |h_B^c = 1|)}$$

Sketch Propagation

- Basic Sketch Propagation
 - Estimate sparsity \hat{s}_C
 - Propagate h_A^r and h_B^c scaled to sparsity \hat{s}_C → E.g., $\sum h_C^r \approx \hat{s}_C ml$
 - Assumption: Structure preserving
- Probabilistic Rounding
 - Ultra-sparse matrices → basic rounding would introduce bias
 - $h_C^r = \text{round}(h_A^r \cdot \hat{s}_C ml / \sum h_A^r)$ w/ probability $h_C^r - [h_C^r]$ of rounding h_C^r up
- Exact Sketch Propagation
 - If A or B are fully diagonal → propagate h_B or h_A , respectively



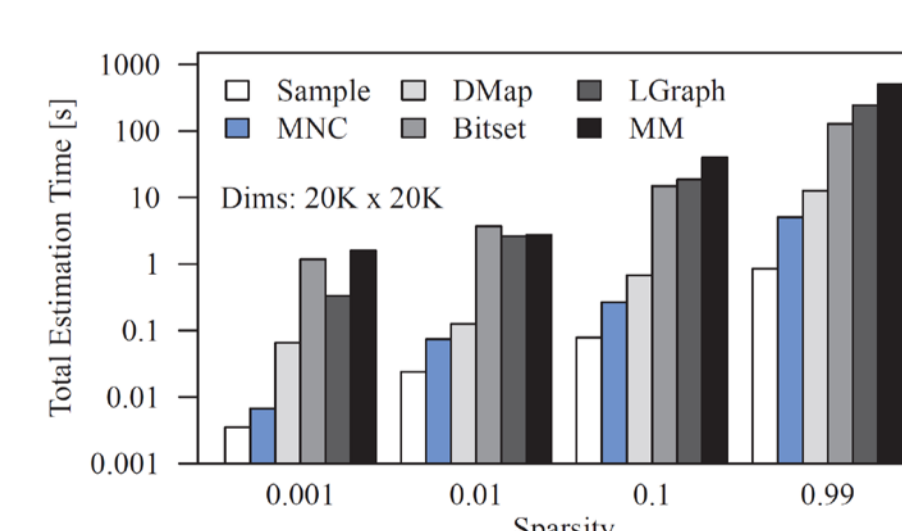
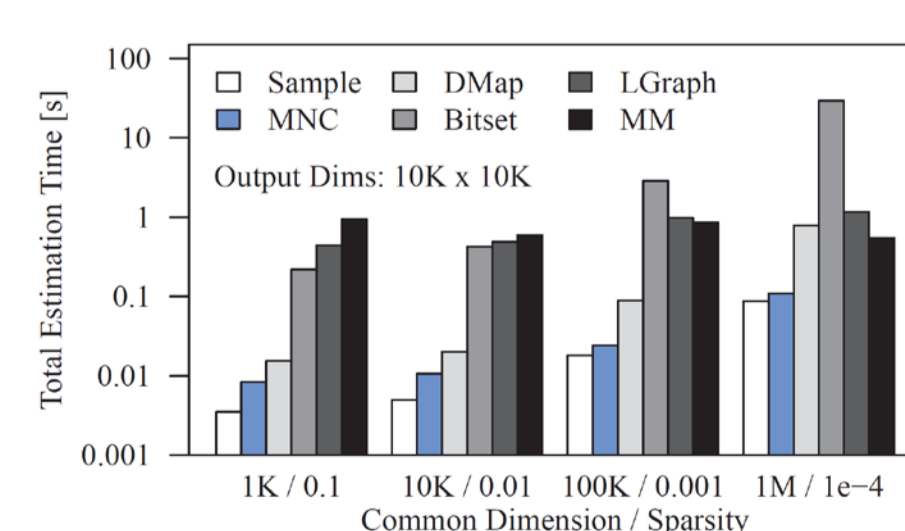
Experiments

SparsEst: Estimation Benchmark

B1 Structured Matrix Products			B2 Real Matrix Operations		
Name	Expression	Data	Name	Expression	Data
B1.1 NLP	X W	Syn1	B2.1 NLP	X W	AMin A
B1.2 Scale	diag(λ) X	Syn2	B2.2 Project	X P	Cov
B1.3 Perm	table(s1,s2) X	Syn2	B2.3 CoRefG	G t(G)	AMin R
B1.4 Outer	C R	Syn3/Syn4	B2.4 EmailG	G G	Email
B1.5 Inner	R C	Syn4/Syn3	B2.5 Mask	M * X	Mnist1m

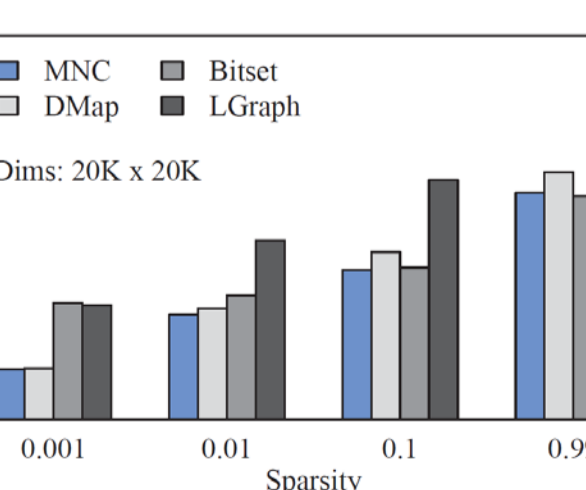
Name	Expression	Data	Dimensions	Sparsity
B3.1 NLP	reshape(X W)	AMin A	25.1M x 2.5M	3.9e-7
B3.2 Scale&Shift	t(S) t(X) diag(w) X S B	Mnist1m	1M x 784	0.25
B3.3 Graph	P G G G G	AMin R	3.1M x 3.1M	2.6e-6
B3.4 Recomm.	(P X != 0) * (P L t(R))	Amazon	8M x 2.3M	1.2e-6
B3.5 Predicate	X * (R * S + T) != 0	Mnist1m		

Runtime Overhead

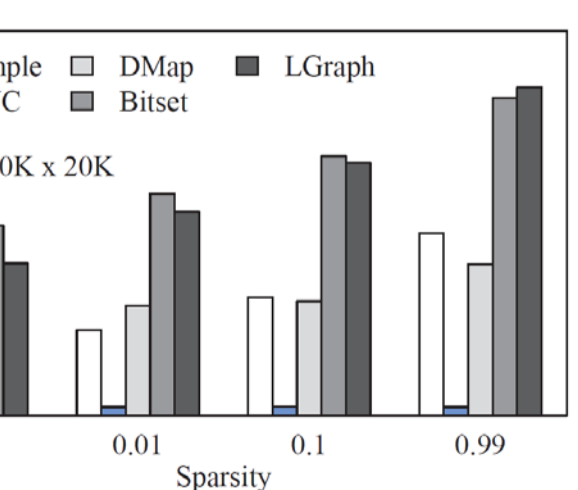


Experimental Setup: 2-socket Xeon E5-2620 (24 vcores, 128GB memory), single-threaded estimators, multi-threaded mm

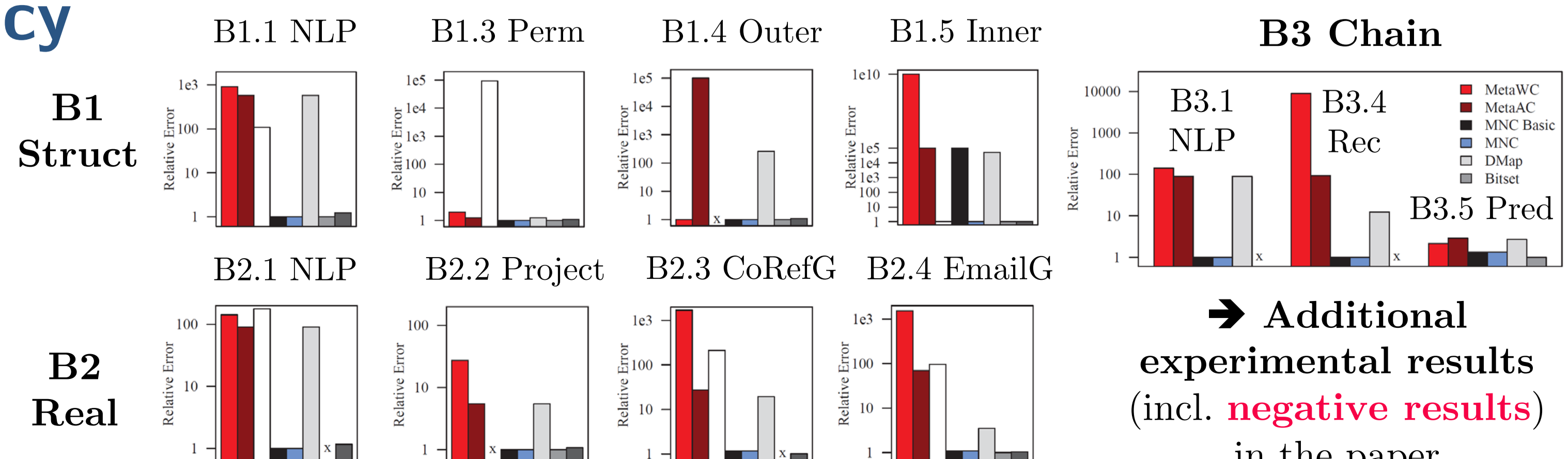
Construction



Estimation



Accuracy



→ Additional experimental results (incl. negative results) in the paper