

Seminar Large-scale Data Engineering (LDE)

01 Structure of Scientific Papers

Dr.-Ing. Patrick Damme

Technische Universität Berlin

Berlin Institute for the Foundations of Learning and Data

Big Data Engineering (DAMS Lab)



Last update: Apr 16, 2023

[Credit: Based on “Introduction to Scientific Writing”/
”01 Structure of Scientific Papers” by Matthias Boehm
(TU Graz, winter 2021/22)]



Announcements/Org



- **Hybrid Setting with Optional Attendance**

- In-person in TEL 811 (~20 seats)
- Virtual via zoom



<https://tu-berlin.zoom.us/j/67376691490?pwd=NmlvWTM5VUVWRjU0UGI2bXhBVkxzQT09>

Agenda



- **FG Big Data Engineering (DAMS Lab)**
- **Course Organization, Outline, and Deliverables**
- **Structure of Scientific Papers**
- **List of Seminar Topics (Proposals)**

FG Big Data Engineering (DAMS Lab)

<https://www.tu.berlin/dams>

FG Big Data Engineering (DAMS Lab) – Team



- **Head of the Group**
Matthias Boehm



- **Postdoc** (01/2021)
Patrick Damme



- **PhD Student** (04/2019)
Arnab Phani



- **PhD Student** (01/2020)
Sebastian Baunsgaard



- **PhD Student** (06/2023)
Christina Dionysio



- **PhD Student** (06/2023)
Philipp Ortner



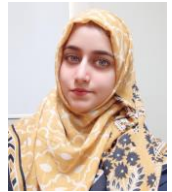
- **PhD Student** (03/2023)
David Justen



- **PhD Student** (02/2023)
Carlos E. Muniz Cuza
[visitor Aalborg University]



- **PhD Student** (09/2019)
Shafaq Siddiqi



- **PhD Student** (04/2021)
Saeed Fathollahzadeh



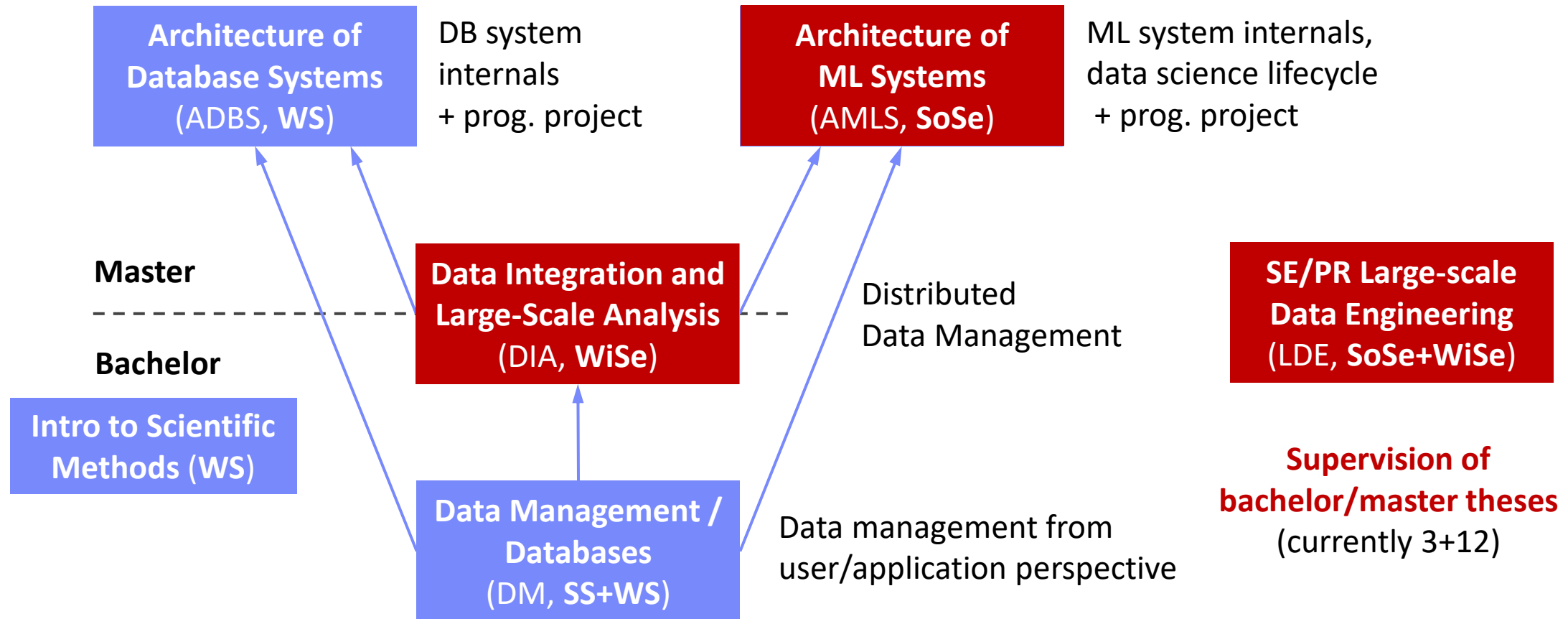
- **3x Student Assistants**
N.N. (1x ~06/2023)

- **Bachelor & Master Students**

About Me

- **Since 10/2022: Postdoc at TU Berlin, Germany**
 - FG Big Data Engineering (DAMS Lab) headed by Prof. Matthias Böhm
 - Continuing work on integrated data analysis pipelines
 - Research interests in the fields of database and ML systems (especially compiler & runtime techniques, extensibility)
- **2021-2022: Postdoc at TU Graz & Know-Center GmbH, Austria**
 - Data Management group headed by Prof. Matthias Böhm
 - Started work on integrated data analysis pipelines
- **2015-2020: PhD student at TU Dresden, Germany**
 - Dresden Database Research Group headed by Prof. Wolfgang Lehner
 - PhD thesis on making complex analytical database queries more efficient through lightweight compression of intermediate results





Course Organization, Outline, and Deliverables

Large-scale Data Engineering: Module Overview



20+5 places in total

bachelor + master

#41086: LDE Seminar + Project (12 ECTS)

13-2 students

12 students

#41095: Seminar LDE (3)

#41094: Project LDE (9 ECTS)

↓
12+2 students

bachelor-only

bachelor-only

Mon, 14:00-16:00
TEL 811 & zoom

Seminar LDE

- Reading & writing scientific papers
- Giving presentations on papers
- Summary paper
- Presentation
- Lecturer & seminar mentor



Project LDE

- Building & evaluating prototypes
- Giving presentations on prototypes
- Prototype design/impl/tests/doc
- Presentation
- Project mentors



Mon, 16:00-18:00
TEL 811 & zoom

- In the context of systems for data engineering, data management, machine learning
- In combination: Ideal preparation for a bachelor/master thesis with our group

Course Organization



■ General Contact Person

- Dr.-Ing. Patrick Damme (patrick.damme@tu-berlin.de)

■ Course Website

- https://pdamme.github.io/teaching/2023_summer/lde/lde_summer2023.html
- One site for seminar and project
- All material, news, **deadlines**, ...

■ Language

- Lectures and slides: **English**
- Communication: **English/German**
- Submitted paper and presentation: **English**
- **Informal language** (first name is fine), immediate feedback is welcome

Semester Schedule & Deadlines



- **Three Introductory Lectures (optional)**
 - Apr 17: Structure of Scientific Papers
 - Apr 24: Scientific Reading and Writing
 - May 8: Experiments, Reproducibility, and Giving Presentations
- **Self-organized Seminar Work**
 - Office hours for any questions (optional)
- **Final Presentations (mandatory)**
 - Jun 26: Session #1
 - Jul 03 : Session #2
 - Jul 10 : Session #3
 - Jul 17 : Session #4
- **List of Seminar Topics**
 - Presented today, take your time to select afterwards
- **Topic Selection**
 - **Deadline: May 1, 23:59 CEST** (in 2 weeks)
 - **First-come-first-serve**
 - List of preferences by email to Patrick Damme
- **Submission of Summary Paper**
 - **Deadline: Jun 19, 23:59 CEST** (in 9 weeks)
 - Summary paper (PDF) by email to Patrick Damme
- **Submission of Presentation Slides**
 - **Deadline: The day before you present, 23:59 CEST**
 - Presentation slides (PDF) by email to Patrick Damme

Seminar Deliverables



▪ Individual Seminar Work

- 1 student = 1 paper, no team work

▪ Summary Paper

- Read and understand selected paper
- Search for related work for some context
- Write summary paper (**6 pages**, excl. references)
 - including related work
 - make sure relation to umbrella is conveyed
- LaTeX with [ACM acmart template](#)
([document class sigconf](#), link on course website)

▪ Presentation

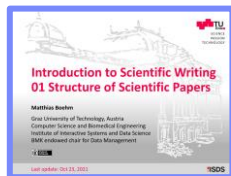
- Summarize your paper
- **15 min talk + 5 min discussion** (stay in time)
- Audience: engage in the discussion

▪ Grading

- Graded portfolio exam
- **#41086 (seminar + project)**
 - 25 pts: summary paper
 - 15 pts: presentation
 - 50 pts: design/impl/tests/doc
 - 10 pts: presentation
- **#41095 (seminar-only)**
 - 65 pts: summary paper
 - 35 pts: presentation

Structure of Scientific Papers

In Computer Science (Data Management)



[**Credit:** Based on “Introduction to Scientific Writing”/
”01 Structure of Scientific Papers” by Matthias Boehm
(TU Graz, winter 2021/22)]

Overview Types of Scientific Writing

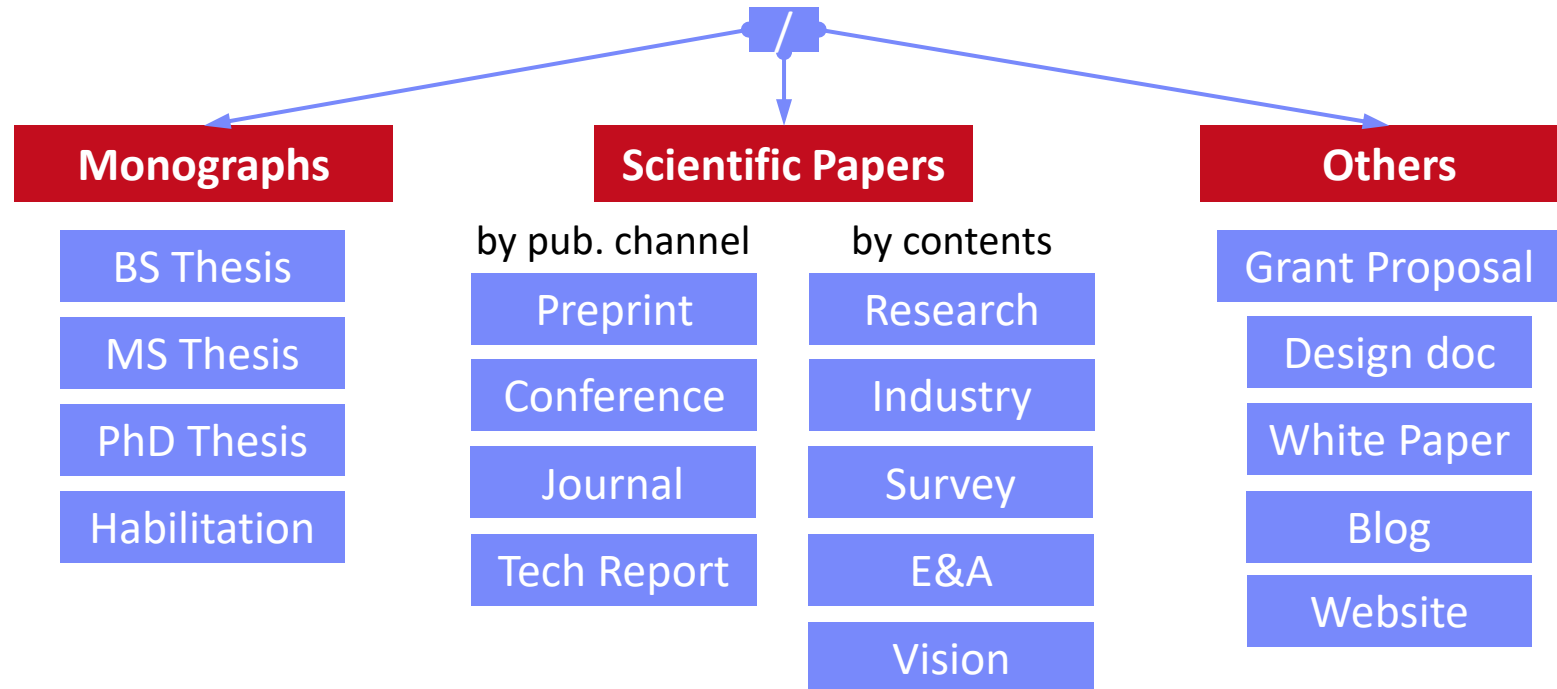


- **Classification of Scientific/Technical Documents**

- Formal vs informal writing, cumulative?, single vs multi-author, archival vs non-archival publications

- **Scientific Writing Skills are crucial**

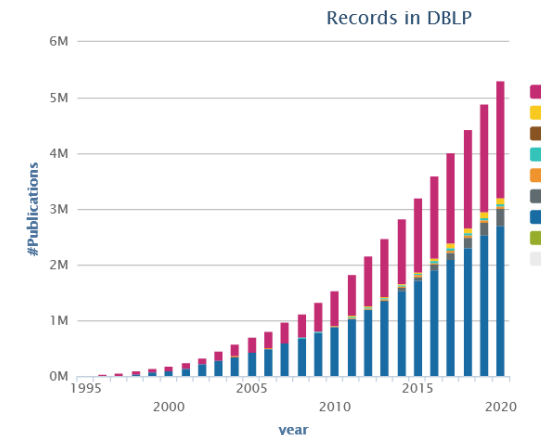
- Different types of docs share many similarities



Preparation



- **Know your Audience**
- **Get your Workflow in Order**
 - **Writing:** LaTeX (e.g., Overleaf, TeXnicCenter), **versioning** (e.g., git), **templates**
 - **Plotting:** R (e.g., plot, ggplot), **Python** (e.g., matplotlib, seaborn), **Gnuplot, LaTeX** (e.g., pgfplots)
 - **Figures:** e.g., MS Visio/MS Powerpoint, **Inkscape** → pdf, eps, svg (vector graphics)
- **Mindset: Quality over Quantity**
 - Aim for top-tier conferences/journals (act as filter)
 - Make the paper useful for others (ideas, evidence, code)



■ Research – Writing Cycle

- Read lots of papers
- ~~Idea~~ → Research → ~~Writing~~ → Document
- Idea → Writing/Research → Document
- Incremental refinement of drafts

■ Paper Submission Cycle

- Blind vs double-blind submission
- Revisions and Camera-ready
- **Similar: bachelor/master** thesis
→ drafts to advisor / final version

■ Example: SIGMOD 2024: Paper Submission Round 2

- **April 15, 2023**: Paper submission
- **May 26 - 28, 2023**: Author feedback phase
- **June 20, 2023**: Notification of accept/reject/review again
- **July 20, 2023**: Revised paper submission
- **August 23, 2023**: Final notification of accept/reject
- **tba, 2023**: Camera ready due

[Recommended Reading]

[Eamonn Keogh: How to do good research, get it published in SIGKDD and get it cited!, **KDD 2009**]



[Simon Peyton Jones: How to write a great research paper, MSR Cambridge]



Compressed Linear Algebra for Large-Scale Machine Learning

Ahmed Elgohary¹, Matthias Boehm¹, Peter J. Haas¹, Frederick R. Reiss¹, Berthold Reinwald²

¹ IBM Research – Almaden; San Jose, CA, USA
² University of Maryland; College Park, MD, USA

ABSTRACT

Large-scale machine learning (ML) algorithms are often iterative, using repeated read-only data access and 1/O-bound matrix-vector multiplications to converge to an optimal model. It is crucial for performance to fit the data into single-node or distributed main memory. General-purpose, heavy- and lightweight compression techniques struggle to achieve both good compression ratios and fast decompression speed to enable block-wise uncompressed operations. Hence, we initiate work on compressed linear algebra (CLA), in which lightweight database compression techniques are applied to matrices and then linear algebra operations such as matrix-vector multiplication are executed directly on the compressed representations. We contribute effective column compression schemes, cache-conscious operations, and an efficient sampling-based compression algorithm. Our experiments show that CLA achieves in-memory operations performance close to the uncompressed case and good compression ratios that allow us to fit larger datasets into available memory. We thereby obtain significant end-to-end performance improvements up to 26x or reduced memory requirements.

1. INTRODUCTION

Data has become a ubiquitous resource [16]. Large-scale machine learning (ML) leverages these large data collections in order to find interesting patterns and build robust predictive models [16, 19]. Applications range from traditional regression analysis and customer classification to recommendations. In this context, often data-parallel frameworks such as MapReduce [20], Spark [5], or Flink [2] are used for cost-effective parallelization on commodity hardware.

Declarative ML: State-of-the-art, large-scale ML aims at declarative ML algorithms [12], expressed in high-level languages, which are often based on linear algebra, i.e., matrix multiplications, aggregations, element-wise and statistical operations. Examples at different abstraction levels are SystemML [21], SciDB [14], Cunnion [27], DMac [30], and TensorFlow [1]. The high level of abstraction gives

*Work done during an internship at IBM Research – Almaden.

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>. For any use beyond those covered by this license, obtain permission by emailing info@bifold.org.

Proceedings of the VLDB Endowment, Vol. 9, No. 12
Copyright 2016 VLDB Endowment 2150-8019/16/08.

960

Example paper used in the following

- Ahmed Elgohary, **Matthias Boehm**, Peter J. Haas, Frederick R. Reiss, Berthold Reinwald: **Compressed Linear Algebra for Large-Scale Machine Learning**. **PVLDB 2016**

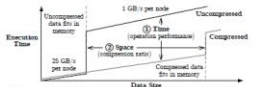


Figure 1: Goals of Compressed Linear Algebra.

data scientists the flexibility to create and customize ML algorithms independent of data and cluster characteristics, without worrying about the underlying data representations (e.g., sparse/dense format) and execution plan generation.

Problem of Memory-Centric Performance: Many ML algorithms are iterative, with repeated read-only access to the data. These algorithms often rely on matrix-vector multiplications to converge to an optimal model. Matrix-vector multiplications are 1/O-bound because they require one complete scan of the matrix, but only two floating point operations per matrix element. Hence, it is crucial for performance to fit the matrix into available memory because memory bandwidth is usually 10x-100x higher than disk bandwidth (but, for matrix-vector, still 10x-40x smaller than peak floating point performance, and thus, matrix-vector remains 1/O-bound). This challenge applies to single-node in-memory computations [28], data-parallel frameworks with distributed caching such as Spark [5], and hardware accelerators like GPUs, with limited device memory [1, 4, 7].

Goals of Compressed Linear Algebra: Declarative ML provides data independence, which allows for automatic compression to fit larger datasets into memory. A baseline solution would be to employ general-purpose compression techniques and decompress matrices block-wise for each operation. However, heavyweight techniques like Crap are not applicable because decompression is too slow, while lightweight methods like Snappy only achieve moderate compression ratios. Existing special-purpose compressed matrix formats with good performance like CSR-V1 [34] similarly show only modest compression ratios. Our approach builds upon research on lightweight database compression, such as compressed bitmaps, and sparse matrix representations. Specifically, we initiate the study of *compressed linear algebra (CLA)*, in which database compression techniques are applied to matrices and then linear algebra operations are executed directly on the compressed representations. Figure 1 shows the goals of this approach: we want to widen the sweet spot for compression by achieving both (1) performance close to uncompressed in-memory operations and (2) good compression ratios to fit larger datasets into memory.

[Ahmed Elgohary, Matthias Boehm, Peter J. Haas, Frederick R. Reiss, Berthold Reinwald: Scaling Machine Learning via Compressed Linear Algebra. **SIGMOD Record 2017 46(1)**]

[Ahmed Elgohary, Matthias Boehm, Peter J. Haas, Frederick R. Reiss, Berthold Reinwald: Compressed Linear Algebra for Large-Scale Machine Learning. **VLDB Journal 2018 27(5)**]

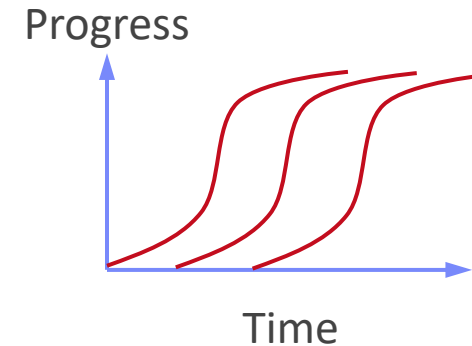
[Ahmed Elgohary, Matthias Boehm, Peter J. Haas, Frederick R. Reiss, Berthold Reinwald: Compressed Linear Algebra for Large-Scale Machine Learning. **Commun. ACM 2019 62(5)**]

Ideas and Topic Selection



■ Problem-Oriented Research

- Focus on problem/observation first, not your solution
- **Discuss early ideas** with collaborators and friends
- Develop your taste for good research topics
- Topic selection needs time → **pipeline model**



■ Ex. Compressed Linear Algebra

- **Problem:** Iterative ML algorithms + memory-bandwidth-bound operations
→ crucial to fit data in memory → automatic lossless compression
- **Sub-problems:** #rows >> #cols, column correlation, column characteristics
→ column-wise compression w/ heterogeneous encoding formats

Prototypes and Experiments

- **Worst Mistake: Schrödinger's Results**
 - Postpone implementation and experiments till last before the deadline
 - No feedback, no reaction time (experiments require many iterations)
 - **Karl Popper:** falsifiability of scientific results
- **Continuous Experiments**
 - Run experiments during survey / prototype building
 - Systematic experiments → observations and ideas for improvements
 - Don't be afraid of throw away prototypes that don't work
- **Ex. Compressed Linear Algebra**
 - Data characteristics inspired overall design of encoding schemes
 - Initially slow compression → dedicated sampling schemes and estimators
 - Initially slow compressed operations → cache-conscious operations, selected operations with better asymptotic behavior

Prototypical Structure of a Scientific Paper



▪ Title & Authors

▪ Sections and Subsections

- Abstract → short overview of problem and solution (part of meta data)
- Introduction → context, problem, contributions
- Background / Preliminaries → necessary background for understanding
- Main Part → your technical core contributions
- Main Part 2
- Experiments → setting, micro benchmarks, end-to-end benchmarks
- Related Work → areas of related work, differences to your own work
- Conclusions → summary, conclusions, and future work
- Acknowledgments → funding agencies, helpful people beyond co-authors
- References → list of other works referenced throughout the paper
- (Appendix) → any additional contents (e.g., proves of theorems, more results)

Title and Authors



List of Authors

- E.g., by contribution (main, ..., advisor)
- E.g., by last name
- Affiliations, contact (corresponding author)

Title

- Descriptive yet concise
- Short name if possible → easier to cite and discuss

Compressed Linear Algebra for Large-Scale Machine Learning

Ahmed Elgohary^{2*}, Matthias Boehm¹, Peter J. Haas¹, Frederick R. Reiss¹,
Berthold Reinwald¹

¹ IBM Research – Almaden; San Jose, CA, USA

² University of Maryland; College Park, MD, USA

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

Tarek Elgamal^{2*}, Shangyu Luo^{3*}, Matthias Boehm¹, Alexandre V. Evfimievski¹,
Shirish Tatikonda⁴, Berthold Reinwald¹, Prithviraj Sen¹

¹ IBM Research – Almaden; San Jose, CA, USA

² University of Illinois; Urbana-Champaign, IL, USA

³ Rice University; Houston, TX, USA

⁴ Target Corporation; Sunnyvale, CA, USA

MNC: Structure-Exploiting Sparsity Estimation for Matrix Expressions

Johanna Sommer
IBM Germany

Matthias Boehm
Graz University of Technology

Alexandre V. Evfimievski
IBM Research – Almaden

Berthold Reinwald
IBM Research – Almaden

Peter J. Haas
UMass Amherst

SliceLine: Fast, Linear-Algebra-based Slice Finding for ML Model Debugging

Svetlana Sagadeeva*
Graz University of Technology

Matthias Boehm
Graz University of Technology

 [Credit: sliceline, Silicon Valley, HBO]

Abstract



% 1. State the problem

Large-scale machine learning (ML) algorithms are often iterative, using repeated read-only data access and I/O-bound matrix-vector multiplications to converge to an optimal model. It is crucial for performance to fit the data into single-node or distributed main memory.

% 2. Say why it's an interesting problem

General-purpose, heavy- and lightweight compression techniques struggle to achieve both good compression ratios and fast decompression speed to enable block-wise uncompressed operations.

% 3. Say what your solution achieves

Hence, we initiate work on compressed linear algebra (CLA), in which lightweight database compression techniques are applied to matrices and then linear algebra operations such as matrix-vector multiplication are executed directly on the compressed representations. We contribute effective column compression schemes, cache-conscious operations, and an efficient sampling-based compression algorithm. Our experiments show that CLA achieves in-memory operations performance close to the uncompressed case and good compression ratios that allow us to fit larger datasets into available memory.

% 4. Say what follows from your solution

We thereby obtain significant end-to-end performance improvements up to 26x or reduced memory requirements.

[Simon Peyton Jones: How to write a great research paper, MSR Cambridge]



ABSTRACT

Large-scale machine learning (ML) algorithms are often iterative, using repeated read-only data access and I/O-bound matrix-vector multiplications to converge to an optimal model. It is crucial for performance to fit the data into single-node or distributed main memory. General-purpose, heavy- and lightweight compression techniques struggle to achieve both good compression ratios and fast decompression speed to enable block-wise uncompressed operations. Hence, we initiate work on compressed linear algebra (CLA), in which lightweight database compression techniques are applied to matrices and then linear algebra operations such as matrix-vector multiplication are executed directly on the compressed representations. We contribute effective column compression schemes, cache-conscious operations, and an efficient sampling-based compression algorithm. Our experiments show that CLA achieves in-memory operations performance close to the uncompressed case and good compression ratios that allow us to fit larger datasets into available memory. We thereby obtain significant end-to-end performance improvements up to 26x or reduced memory requirements.

Introduction



■ Prototypical Structure

- Context (1 paragraph)
- Problems (1-3 paragraphs)
- [Existing Work (1 paragraph)]
- [Idea (1 paragraph)]
- Contributions (1 paragraph)



■ Introduction Matters

- **Anchoring:** most reviewers reach their opinion after reading introduction and motivation and then look for evidence

Contributions: Our major contribution is to make a case for *compressed linear algebra*, where linear algebra operations are directly executed over compressed matrices. We leverage ideas from database compression techniques and sparse matrix representations. The novelty of our approach is a combination of both, leading towards a generalization of sparse matrix representations and operations. The structure of the paper reflects our detailed technical contributions:

- **Workload Characterization:** We provide the background and motivation for CLA in **Section 2** by giving an overview of Apache SystemML, and describing typical linear algebra operations and data characteristics.
- **Compression Schemes:** We adapt several column-based compression schemes to numeric matrices in **Section 3** and describe efficient, cache-conscious core linear algebra operations over compressed matrices.
- **Compression Planning:** In **Section 4**, we further provide an efficient sampling-based algorithm for selecting a good compression plan, including techniques for compressed-size estimation and column grouping.
- **Experiments:** Finally, we integrated CLA into Apache SystemML. In **Section 5**, we study a variety of full-fledged ML algorithms and real-world datasets in both single-node and distributed settings. We also compare CLA against alternative compression schemes.

[Eamonn Keogh: How to do good research, get it published in SIGKDD and get it cited!, KDD 2009]



Writing the Paper (and more Experiments)



■ Easily Readable: Quality \propto Time

■ Make it easy to skim the paper

→ paragraph labels, self-explanatory figures (close to text), and structure

- Avoid unnecessary formalism → as simple as possible
- Shortening the text in favor of structure improves readability

■ Ex. Compressed Linear Algebra

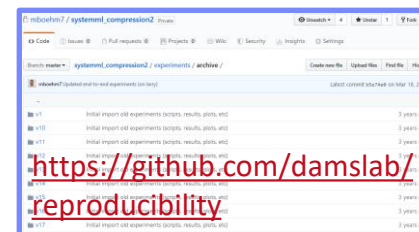
- Initial SIGMOD submission: **12+3 pages**
- Final PVLDB submission: **12 pages**
(+ more figures, experiments, etc.)



→ 02 Scientific Reading and Writing

■ Solid, Reproducible Experiments

- Create, use, and share dedicated benchmarks / datasets
- Avoid weak baselines, start early w/ baseline comparisons
- Automate your experiments as much as possible
- Keep repository of all scripts, results, and used parameters



→ 03 Experiments, Reproducibility, and Giving Presentations

Related Work



■ Purpose of a “Related Work”-Section

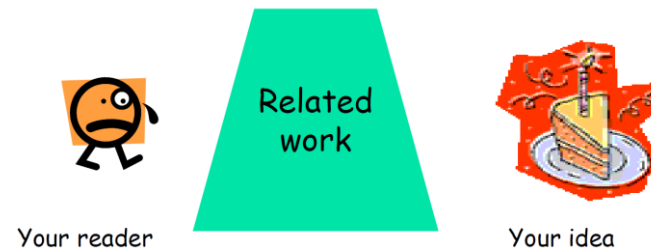
- **Not** a mandatory task or to show you know the field
- Put your work in context of related areas (~ 1 paragraph each)
- Discuss closely related work
- **Crisp separation from existing work** (what are the differences)

[Simon Peyton Jones: How to write a great research paper, MSR Cambridge]



■ Placement

- Section 2 or **Section n-1**
- Throughout the paper



■ Give Credit

- Cite broadly, **give credit to inspiring ideas**, create connections
- Honestly acknowledge **limitations of your approach**

References

Setup

- Use LaTeX `\cite{}` and BibTeX
- Use a consistent source of bibtex entries (e.g., DBLP)

```
inproceedings{StonebrakerBPR11,
  author    = {Michael {Stonebraker et al.}},
  title     = {{The Architecture of SciDB}},
  booktitle = {{SSDBM}},
  year      = {2011}
}
```

VLDB2016.bib

`\bibliographystyle{abbrv}`
`\bibliography{VLDB2016}`

Different References Styles

- But, **not in footnotes** (unless required)

8. REFERENCES

- M. Abadi et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. *CoRR*, 2016.
- A. Alexandrov et al. The Stratosphere Platform for Big Data Analytics. *VLDB J.*, 23(6), 2014.

References

- [All14] Alexandrov, A. et al.: The Stratosphere platform for big data J. 23/6, 2014.
- [AS14] Arap, O.; Swany, M.: Offloading MPI Parallel Prefix Scan the NetFPGA. CoRR abs/1408.4939/, 2014.

7. CONCLUSIONS

We have initiated work on compressed linear algebra (CLA), in which matrices are compressed with lightweight techniques and linear algebra operations are performed directly on the compressed representation. We introduced effective column encoding schemes, efficient operations on compressed matrices, and an efficient sampling-based compression algorithm. Our experiments show operations performance close to the uncompressed case and compression ratios similar to heavyweight formats like Cray but better than lightweight formats like Snappy, providing significant performance benefits when data does not fit into memory. Thus, we have demonstrated the general feasibility of CLA, enabled by declarative ML that hides the underlying physical data representation. CLA generalizes sparse matrix representations, encoding both dense and sparse matrices in a universal compressed form. CLA is also broadly applicable to any system that provides blocked matrix representations, linear algebra, and physical data independence. Interesting future work includes (1) full optimizer integration, (2) global planning and physical design tuning, (3) alternative compression schemes, and (4) operations beyond matrix-vector.

8. REFERENCES

- M. Abadi et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. *CoRR*, 2016.
- A. Alexandrov et al. The Stratosphere Platform for Big Data Analytics. *VLDB J.*, 23(6), 2014.
- A. Aghajani et al. An Efficient Two-Dimensional Blocking Strategy for Sparse Matrix-Vector Multiplication on GPUs. In *ICS (Int. Conf. on Supercomputing)*, 2014.
- A. Aghajani et al. On Optimizing Machine Learning Workloads via Kernel Fusion. In *PPoPP (Principles and Practice of Parallel Programming)*, 2015.
- M. A. Bostrom. Data Compression in Scientific and Statistical Databases. *TSE (Trans. SIP Eng.)*, 11(10), 1985.
- N. Bell and M. Garland. Implementing Sparse Matrix-Vector Multiplication on Throughput-Oriented Processors. In *SC (Supercomputing Conf.)*, 2009.
- J. Bergstra et al. Theano: a CPU and GPU Math Expression Compiler. In *SciPy*, 2010.
- K. S. Boyer et al. On Synopses for Distinct-Value Estimation Under Matrix Operations. In *SIGMOD*, 2007.
- B. Bhattacharjee et al. Efficient Index Compression in DR2 L3W. *PVLDB*, 2(2), 2009.
- S. Bhattacharjee et al. PISos: An Efficient Storage Framework for Managing Scientific Data. In *SEDDM*, 2014.
- C. Huang et al. Dictionary-based Order-preserving String Compression for Main Memory Column Stores. In *SIGMOD*, 2009.
- M. Bilenko et al. Declarative Machine Learning - A Classification of Basic Properties and Types. *CoRR*, 2016.
- L. Boston. The infinite MNIST dataset. <http://lambertson.com/projects/infmnist>.
- M. Charlier et al. Towards Estimation Error Guarantees for Distinct Values. In *SIGMOD*, 2000.
- R. Chittka et al. Approximate Kernel k-means: Solution to Large Scale Kernel Clustering. In *KDD*, 2011.
- J. Cohen et al. MAD Skills: New Analysis Practices for Big Data. *PVLDB*, 2(2), 2009.
- C. Constantinides and M. Lu. Quick Estimation of Data Compression and De-duplication for Large Storage Systems. In *CCP (Data Compression, Clust. and Proc.)*, 2011.
- C. V. Curmeck. Data Compression on a Database System. *Commun. ACM*, 28(12), 1985.
- S. Das et al. Riscv: Integrating R and Hadoop. In *SIGMOD*, 2010.
- J. Dean and S. Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. In *OSDI*, 2004.

- A. Ghosh et al. SystemML: Declarative Machine Learning on MapReduce. In *ICDE*, 2011.
- J. J. Good. The Population Frequency of Species and the Estimation of Population Parameters. *Biometrika*, 1953.
- G. Graefe and L. D. Shapiro. Data Compression and Database Performance. In *Applied Computing*, 1991.
- P. J. Haas and L. S. Stokes. Estimating the Number of Clones in a Finite Population. *J. Amer. Statist. Assoc.*, 98(444), 1998.
- D. Harsh et al. Estimation of Deduplication Ratios in Large Data Sets. In *MSST (Mass Storage Sys. Tech.)*, 2012.
- D. Harsh et al. To Zip or not to Zip: Effective Resource Usage for Real-Time Compression. In *PAST*, 2013.
- B. Huang et al. Cambda: Optimizing Statistical Data Analysis in the Cloud. In *SIGMOD*, 2013.
- B. Huang et al. Resource Elasticity for Large-Scale Machine Learning. In *SIGMOD*, 2015.
- S. Ilieva et al. Estimating the Compression Fraction of an Index using Sampling. In *ICDE*, 2010.
- K. L. Johnson et al. *Transformed Discrete Distributions*. Wiley, New York, 2nd edition, 1992.
- V. Kulkarni et al. An Extended Compression Format for the Optimization of Sparse Matrix-Vector Multiplication. *TPDS (Trans. Par. and Dist. Systems)*, 24(10), 2013.
- D. Korrer et al. SLACE: Sparse Linear Algebra in a Column-Oriented In-Memory Database System. In *SEDDM*, 2014.
- H. Kimura et al. Compression Aware Physical Database Design. *PVLDB*, 4(10), 2011.
- K. Kouris et al. Optimizing Sparse Matrix-Vector Multiplication Using Index and Value Compression. In *CF (Computing Frontiers)*, 2008.
- H. Lang et al. Data Blocks: Hybrid OLTP and OLAP on Compressed Storage using both Vectorization and Compression. In *SIGMOD*, 2016.
- P. Larson et al. SQL Server Column Store Indexes. In *SIGMOD*, 2011.
- M. Leisner. UCI Machine Learning Repository. Higgs, Coorsys, CS Course (1990). archive.ics.uci.edu/ml/.
- P. E. O'Neill. Model 204 Architecture and Performance. In *High Performance Transaction Systems*, 1980.
- Oracle. *Data Warehousing Guide, 11g Release 1*, 2007.
- V. Raman and G. Swart. How to Write a Table Driven Compressed Relation. In *VLDB*, 2006.
- V. Raman et al. DR2 with REL Acceleration: So Much More than Just a Column Store. *PVLDB*, 6(11), 2013.
- Y. Saito. SPARSKIT: a basic tool kit for sparse matrix computations - Version 2, 1994.
- M. Stonebraker et al. C-Store: A Column-oriented DBMS. In *VLDB*, 2005.
- M. Stonebraker et al. The Architecture of SciDB. In *SEDDM*, 2011.
- System. *IQ 15.4 System Administration Guide*, 2013.
- G. Valsani and F. Valenti. Estimating the Unseen: An ϵ -Neyman-Savage Estimator for Entropy and Support Size, Shows Optimal via New CLTs. In *STOC*, 2011.
- F. Vassiliadis et al. The Implementation and Performance of Compressed Databases. *SIGMOD Record*, 29(3), 2000.
- S. Williams et al. Optimization of Sparse Matrix-Vector Multiplication on Emerging Multiscale Platforms. In *SC (Supercomputing Conf.)*, 2007.
- K. Wu et al. Optimizing Bitmap Indices With Efficient Compression. *TPDS*, 31(1), 2006.
- L. Ye et al. Exploring Matrix Dependency for Efficient Distributed Matrix Computation. In *SIGMOD*, 2015.
- M. Zaharia et al. Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing. In *NIPS*, 2012.
- C. Zhang et al. Materialized Optimization for Feature Selection Workloads. In *SIGMOD*, 2014.

971

References

Jaume Amores. 2013. Multiple instance classification: Review, taxonomy and comparative study. *Artificial Intelligence*.



Dealing with Feedback / Criticism



■ Different Kinds of Feedback

- Casual discussion of early ideas
- Comments on paper drafts
- Reviewer comments (good and bad)

- Always welcome feedback/criticism
- Address all feedback w/ sincere effort

■ Example Compressed Linear Algebra

- SIGMOD Reviewer 2 (REJECT)

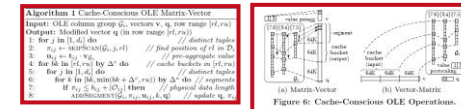
- “The rewriting for $q=Xv$ seems wrong: To compute q , one takes each row of the matrix X and multiplies it with the vector v .”

- PVLDB Reviewer 3 (WEAK ACCEPT)

- “I kinda disagree with the broad definition of declarative ML [...]”

■ Paper Rebuttal and/or Revision

- **Rebuttal**: seriously consider all feedback (in doubt agree), and answer with facts / ideas how to address the comments
- **Revision** (conditional accept): address all revision requests



large n . We therefore use cache-conscious schemes for OLE and RLE based on horizontal, segment-aligned scans (with benefits of up to 2.1x/5.4x for M/V/V/M in our experiments) see Algorithm 1 and Figure 6(a) for the case of OLE. Multi-threaded operations parallelize over segment-aligned partitions of rows $[r, r+a)$, which guarantees disjoint results and thus avoids partial results per thread. We find v_{i_1} , the starting position of each v_{i_2} in D , via a skip scan that aggregates segment lengths until we reach i (line 2). To minimize the overhead of finding v_{i_1} , we use static scheduling (task partitioning). We further precompute $h_{i_1} = h_{i_2} + v_{i_1}$ once for all tuples (line 3). For each cache-bucket of size Δ^* (such that Δ^* divides Δ in \mathbb{N} and by default $\Delta^* = 2\Delta^*$), we then iterate over all distinct tuples (line 5-8) but maintain the current position n , as well. The inner non-zero p_i is written to both R_{i_1} and R_{i_2} . Multi-threaded operations dynamically parallelize over column groups.

Other Operations: Various common operations can be executed very efficiently over compressed matrices without scanning the offset lists. Sparse-side matrix-scalar operations such as $X^T \cdot \alpha$ or αX are carried out with a single pass over the set of right T_i for each column group G_i . Aggregations such as adding a column of 1's or another matrix to X is done via simple concatenation of column groups. Finally, unary aggregates like sum (or similarly colSums) are efficiently computed via routines like $\sum_{i=1}^n \sum_{j=1}^m \text{CO}(X)_{ij}$. For each value, we aggregate the RLE run lengths or OLE lengths per segment, respectively. Row aggregates (e.g., rowSums) are computed in a cache-conscious manner.



[Matthias Boehm et al.: Declarative Machine Learning - A Classification of Basic Properties and Types. **CoRR 2016.**]



List of Seminar Topics (Proposals)

Last update: Apr 16, 2023

- Papers in the fields of **systems** for **data engineering**, **data management**, and **machine learning**
- This Semester's **Umbrella Topic: "Extensible Data Systems"**
 - Motivation
 - Meet requirements of **emerging applications** (multimedia, machine learning, domain-specific, ...)
 - Enable the timely **adoption of novel techniques and technologies** (algorithms, hardware, ...)
 - Handle the **increasing specialization** (data models, data representations, algorithms, hardware, ...)
 - **Facilitate research** at systems and algorithms level
 - Active field of research for decades (at least 1980s till today)
 - Interesting challenges
 - **Which abstractions** are needed?
 - **How to enable users** to extend the system?
 - **How to reach efficiency** in spite of these abstractions?
 - Concepts have been proposed at all levels of the system stack

▪ Extensible Database Management Systems

- Implementation of Data Abstraction in the Relational Database System Ingres (SIGMOD Rec, 1984)
- The Design of POSTGRES (SIGMOD, 1986)
- Design and Implementation of an Extensible Database Management System Supporting User Defined Data Types and Functions (VLDB, 1988)
- Starburst Mid-Flight: As the Dust Clears (IEEE Trans. Knowl. Data Eng., 1990)
- MIL Primitives for Querying a Fragmented World (VLDBJ, 1999)
- Hyrise Re-engineered: An Extensible Database System for Research in Relational In-Memory Data Management (EDBT, 2019)

▪ Extensible Machine Learning Systems

- TensorFlow: A System for Large-Scale Machine Learning (OSDI, 2016)
- BAGUA: Scaling up Distributed Learning with System Relaxations (PVLDB, 2022)

- **Component-based Database/ML Systems**
 - Flashlight: Enabling Innovation in Tools for Machine Learning (PMLR, 2022)
 - mutable: A Modern DBMS for Research and Fast Prototyping (CIDR, 2023)
- **Optimizer Generators / Compiler Frameworks**
 - GENESIS: An Extensible Database Management System (IEEE Trans. Software Eng., 1988)
 - The Volcano Optimizer Generator: Extensibility and Efficient Search (ICDE, 1993)
 - MLIR: Scaling Compiler Infrastructure for Domain Specific Computation (CGO, 2021)
- **Efficient Handling of User-defined Functions**
 - Opening the Black Boxes in Data Flow Optimization (PVLDB, 2012)
 - Extending database task schedulers for multi-threaded application code (SSDBM, 2015)
 - Processing Java UDFs in a C++ Environment (SOCC, 2017)

Seminar Topics (3/4)



▪ Extensibility of Hardware Backends

- Voodoo - A Vector Algebra for Portable Database Performance on Modern Hardware (PVLDB, 2016)
- TVM: An Automated End-to-End Optimizing Compiler for Deep Learning (OSDI, 2018)
- Hardware-Oblivious SIMD Parallelism for In-Memory Column-Stores (CIDR, 2020)

▪ Extensible Data Analytics

- Towards a unified architecture for in-RDBMS analytics (SIGMOD, 2012)
- MLbase: A Distributed Machine-learning System (CIDR, 2013)

Disclaimer:
Co-authored by
the lecturer

Seminar Topics (4/4)



- **Miscellaneous**

- MorphStore: Analytical Query Engine with a Holistic Compression-Enabled Processing Model (PVLDB, 2020)
- User-Defined Operators: Efficiently Integrating Custom Algorithms into Modern Databases (PVLDB, 2022)

Disclaimer:
Co-authored by
the lecturer

- **Alternative: Propose yet another topic/paper**

- Should be related to the umbrella topic of extensible data systems

Summary and Q&A



- **FG Big Data Engineering (DAMS Lab)**
- **Course Organization, Outline, and Deliverables**
- **Structure of Scientific Papers**
- **List of Seminar Topics (Proposals)**

- **Remaining Questions?**

- **Next Lectures**
 - 02 **Scientific Reading and Writing** [Apr 24]
 - 03 **Experiments, Reproducibility, and Giving Presentations** [May 8]