

# 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 14, 2024

[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

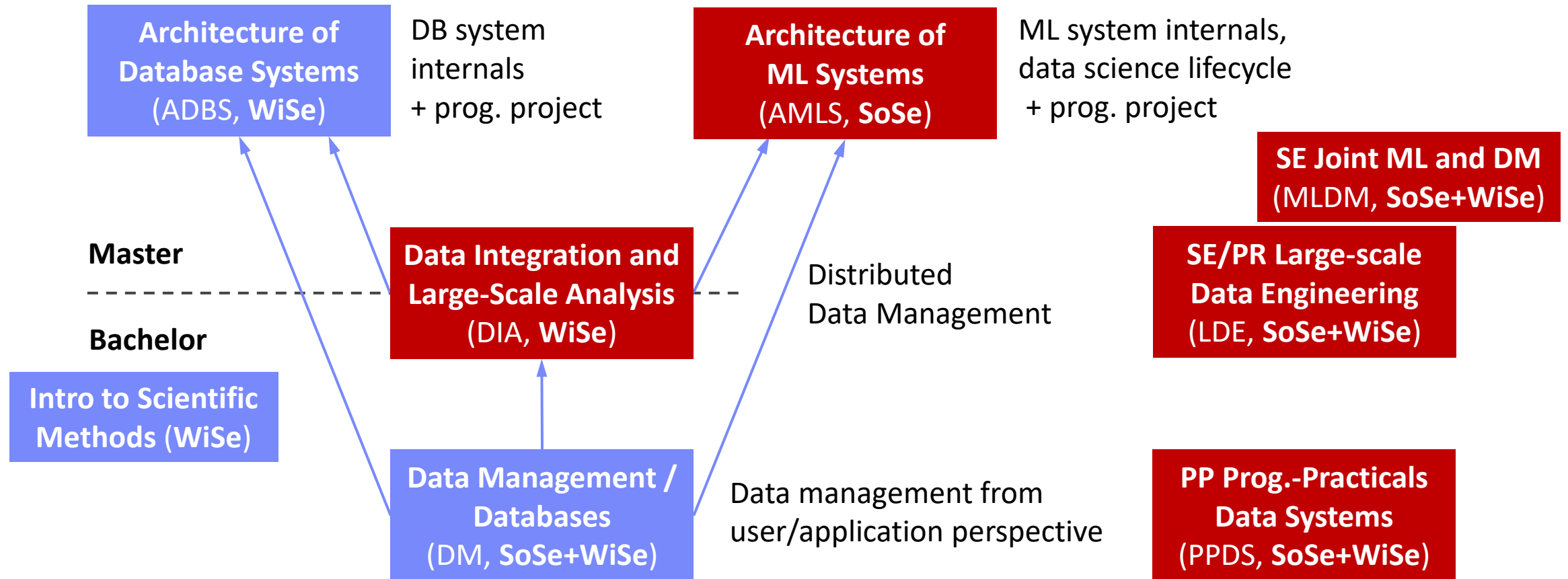


<https://tu-berlin.zoom.us/j/67376691490?pwd=NmlvWTM5VUVWRjU0UGI2bXhBVkxzQT09>

# 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





# Agenda



- **Course Organization, Outline, and Deliverables**
- **Structure of Scientific Papers**
- **List of Seminar Topics**

# Course Organization, Outline, and Deliverables

# Large-scale Data Engineering: Module Overview



20 places in total

bachelor + master

#41086: LDE Seminar + Project (12 ECTS)

10 students

10 students

#41095: Seminar LDE (3)

#41183: Project LDE (9 ECTS)

8 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](mailto:patrick.damme@tu-berlin.de))

## ■ Course Website

- [https://pdamme.github.io/teaching/2024\\_summer/lde/lde\\_summer2024.html](https://pdamme.github.io/teaching/2024_summer/lde/lde_summer2024.html)
- One site for seminar and project
- All material, schedule, **deadlines**

## ■ ISIS course

- <https://isis.tu-berlin.de/course/view.php?id=37443>
- Announcements, discussion forum, polls for topic selection

## ■ 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 15: Structure of Scientific Papers
  - Apr 22: *no lecture, office hour in TEL 814*
  - Apr 29: Scientific Reading and Writing
  - May 06: Experiments, Reproducibility, and Giving Presentations
- **Self-organized Seminar Work**
  - Office hours for any questions (optional)
- **Final Presentations (mandatory)**
  - Jul 01, 14:00-18:00: Session #1
  - Jul 15, 14:00-18:00: Session #2
- **List of Seminar Topics**
  - Presented today, take your time to select afterwards
- **Topic Selection**
  - **Deadline: Apr 29, 23:59 CET** (in 2 weeks)
  - Ranked list of **5 topics** via poll on the ISIS course
  - Global topic assignment based on preferences
  - **Notification of assigned topics: May 06** (in 3 weeks)
- **Submission of Summary Paper**
  - **Deadline: Jun 24, 23:59 CET** (in 10 weeks)
  - Summary paper (PDF) by email to Patrick Damme
- **Submission of Presentation Slides**
  - **Deadline: The day before you present, 23:59 CET**
  - Presentation slides (PDF) by email to Patrick Damme

# Seminar Deliverables



- **Individual Seminar Work**
  - 1 student = 1 paper, no team work
- **Summary Paper (in English)**
  - Read and understand selected paper
  - Search for related work to provide some context
  - Write summary paper (**4 pages**, excl. references)
    - including related work
    - make sure relation to umbrella topic is conveyed
  - LaTeX with given template
- **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
- **Academic Honesty / No Plagiarism**  
(incl LLMs like ChatGPT)

## ACM acmart template document class sigconf (double-column)

**CCS concepts**  
(ACM Computing Classification System)  
Select concepts at <http://dl.acm.org/ccs.cfm>  
and insert generated code:

```

220 %\
221 %% The code below is generated by the tool at http://dl.acm.org/ccs.cfm.
222 %% Please copy and paste the code instead of the example below.
223 %%
224 \begin{CCSXML}
225 <ccs2012>
226 <concept>
227 <concept_id>10010520.10010553.10010562</concept_id>
228 <concept_desc>Computer systems organization-Embedded systems</concept_desc>
229 <concept_significance>500</concept_significance>
230 </concept>
231 <concept>
232 <concept_id>10010520.10010575.10010755</concept_id>
233 <concept_desc>Computer systems organization-Redundancy</concept_desc>
234 <concept_significance>300</concept_significance>
235 </concept>
236 </CCSXML>
    
```

**Copyright notice**  
Just keep as it is (ignore the dummy data)

### The Name of the Title Is Hope

Ben Trovato* G.K.M. Tobin* trovato@corporation.com webmaster@marysville-ohio.com Institute for Clarity in Documentation Dublin, Ohio, USA	Lars Thorvöld The Thorvöld Group Hekla, Iceland larst@affiliation.org	Valerie Béranger Inria Paris-Rocquencourt Rocquencourt, France
Aparna Patel Rajiv Gandhi University Doimukh, Arunachal Pradesh, India	Huifen Chan Tsinghua University Haidian Qu, Beijing Shi, China	Charles Palmer Palmer Research Laboratories San Antonio, Texas, USA cpalmer@prl.com
John Smith The Thorvöld Group Hekla, Iceland jsmith@affiliation.org	Julius P. Kumquat The Kumquat Consortium New York, USA jpkumquat@consortium.net	




Figure 1: Seattle Mariners at Spring Training, 2010.

**ABSTRACT**  
A clear and well-documented L<sup>A</sup>T<sub>E</sub>X document is presented as an article formatted for publication by ACM in a conference proceedings or journal publication. Based on the "acmart" document class, this article presents and explains many of the common variations, as well as many of the formatting elements an author may use in the preparation of the documentation of their work.

**CCS CONCEPTS**  
• Computer systems organization → Embedded systems; Redundancy; Robotics; • Networks → Network reliability.

\*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.  
Conference acronym 'XX, June 01–05, 2018, Woodstock, NY  
© 2018 Association for Computing Machinery  
ACM ISBN 978-1-60959-XXXX-X/18/06... \$15.00  
<https://doi.org/XXXXXX.XXXXXX>

**KEYWORDS**  
datasets, neural networks, gaze detection, text tagging

**ACM Reference Format:**  
Ben Trovato, G.K.M. Tobin, Lars Thorvöld, Valerie Béranger, Aparna Patel, Huifen Chan, Charles Palmer, John Smith, and Julius P. Kumquat. 2018. The Name of the Title Is Hope. In *Proceedings of the conference on the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 6 pages. <https://doi.org/XXXXXX.XXXXXX>

**1 INTRODUCTION**  
ACM's consolidated article template, introduced in 2017, provides a consistent L<sup>A</sup>T<sub>E</sub>X style for use across ACM publications, and incorporates accessibility and metadata-extraction functionality necessary for future Digital Library endeavors. Numerous ACM and SIG-specific L<sup>A</sup>T<sub>E</sub>X templates have been examined, and their unique features incorporated into this single new template.

If you are new to publishing with ACM, this document is a valuable guide to the process of preparing your work for publication. If you have published with ACM before, this document provides insight and instruction into more recent changes to the article template.

<https://www.acm.org/publications/proceedings-template>

**Teaser image**  
Not required (especially no photograph)

**Keywords**  
Specify meaningful keywords

```

255 %% Keywords. The author(s) should pick words that accurately describe
256 %% the work being presented. Separate the keywords with commas.
257 \keywords{datasets, neural networks, gaze detection, text tagging}
    
```

**ACM reference format**  
Just keep as it is (ignore the dummy data)

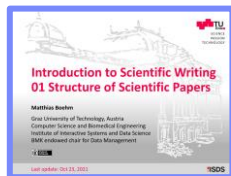
# Portfolio Exam Registration



- **Portfolio exam registration: May 06 – Jun 03**
  - Binding registration in Moses/MTS
  - Including selection of seminar presentation date (first-come-first-serve)
- **Portfolio exam de-registration**
  - **Until 3 days before the first graded exam part**
    - Modules “LDE”/”Seminar LDE”: until **Jun 21**
    - Module “Project LDE”: until **Jul 26**
    - De-register yourself in Moses/MTS
  - **With sufficient reason: Until the day of the exam**
    - In case of sickness etc.
    - Modules “LDE”/”Seminar LDE”: until **Jun 24**
    - Module “Project LDE”: until **Jul 29**
- **Missing deadlines/exam without de-registration**
  - Zero points in the respective exam part (!)
  - **Approach us early in case of problems**
- **If you don't want to take LDE anymore**
  - Let me know asap to give students in the queue a chance to fill in

# 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

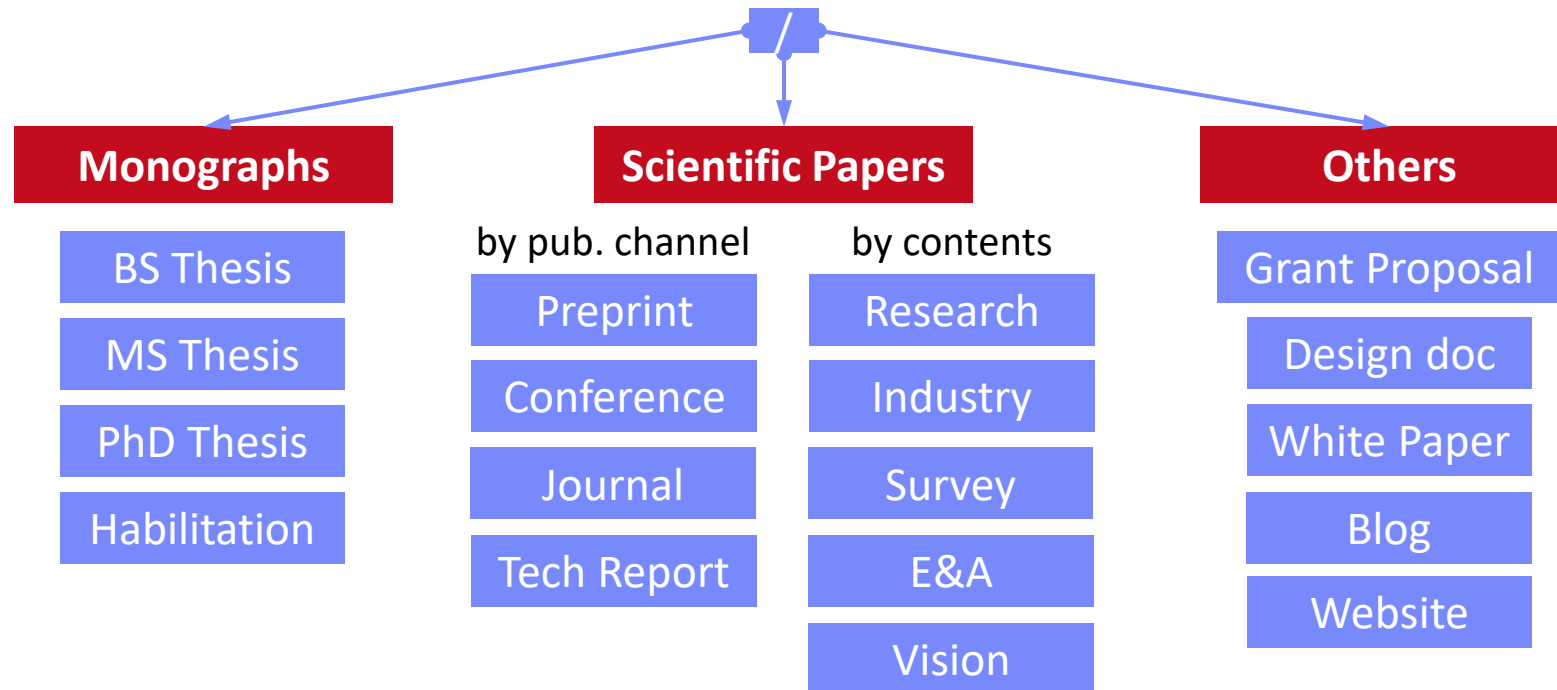


- **Classification of Scientific/Technical Documents**

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

- **Scientific Writing Skills are crucial**

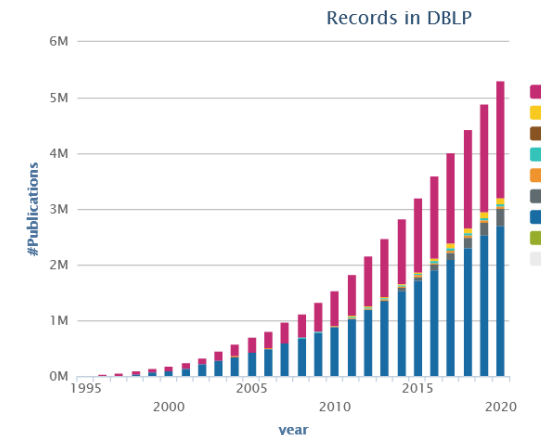
- Different types of docs share many similarities



# 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 Elgohary<sup>1</sup>, Matthias Boehm<sup>1</sup>, Peter J. Haas<sup>1</sup>, Frederick R. Reiss<sup>1</sup>, Berthold Reinwald<sup>2</sup>

<sup>1</sup> IBM Research – Almaden; San Jose, CA, USA  
<sup>2</sup> 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](mailto:info@bifold.org).

Proceedings of the VLDB Endowment, Vol. 9, No. 12  
Copyright 2016 VLDB Endowment 2150-8099/16/08.

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**



[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)**]

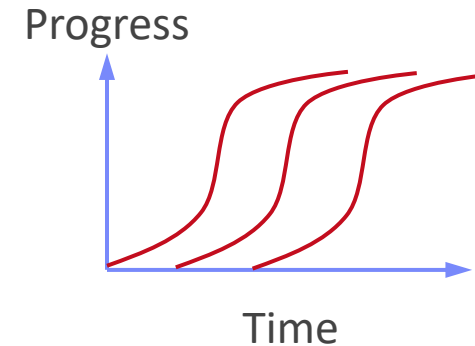
[Ahmed Elgohary, Matthias Boehm, Peter J. Haas, Frederick R. Reiss, Berthold Reinwald: Compressed Linear Algebra for Large-Scale Machine Learning. **Commun. ACM 2019 62(5)**]

# 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 throwing 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<sup>2\*</sup>, Matthias Boehm<sup>1</sup>, Peter J. Haas<sup>1</sup>, Frederick R. Reiss<sup>1</sup>,  
Berthold Reinwald<sup>1</sup>

<sup>1</sup> IBM Research – Almaden; San Jose, CA, USA

<sup>2</sup> University of Maryland; College Park, MD, USA

## 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; San Jose, CA, USA

<sup>2</sup> University of Illinois; Urbana-Champaign, IL, USA

<sup>3</sup> Rice University; Houston, TX, USA

<sup>4</sup> 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)



**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.

## ■ Introduction Matters

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

[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

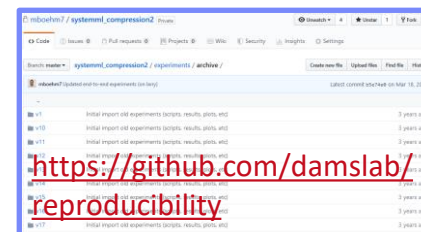
- Initial SIGMOD submission: **12+3 pages**
- Final PVLDB submission: **12 pages**  
(+ more figures, experiments, etc.)



→ 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

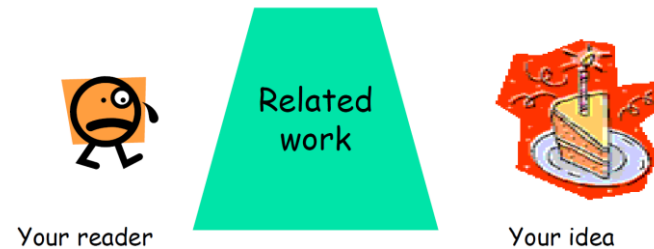
- **Not** a mandatory task or to show you know the field
- Put your work in context of related areas (~ 1 paragraph each)
- Discuss closely related work
- **Crisp separation from existing work** (what are the differences)

[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

[1] M. Abadi et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. *CoRR*, 2016.

[2] A. Alexandrov et al. The Stratosphere Platform for Big Data Analytics. *VLDB J.*, 23(6), 2014.

### References

[All4] 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 over 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

[1] M. Abadi et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. *CoRR*, 2016.

[2] A. Alexandrov et al. The Stratosphere Platform for Big Data Analytics. *VLDB J.*, 23(6), 2014.

[3] A. Aghajari et al. An Efficient Two-Dimensional Blocking Strategy for Sparse Matrix-Vector Multiplication on GPUs. In *ICS (Int. Conf. on Supercomputing)*, 2014.

[4] A. Aghajari et al. On Optimizing Machine Learning Workloads via Kernel Fusion. In *PPoPP (Principles and Practice of Parallel Programming)*, 2015.

[5] M. A. Broomhead. Data Compression in Scientific and Statistical Databases. *TSE Trans. SIV Eng J.*, 11(10), 1985.

[6] N. Bell and M. Garland. Implementing Sparse Matrix-Vector Multiplication on Throughput-Oriented Processors. In *SC (Supercomputing Conf.)*, 2009.

[7] J. Bergstra et al. Theano: a CPU and GPU Math Expression Compiler. In *SciPy*, 2010.

[8] K. S. Boyer et al. On Synapses for Distinct-Value Estimation Under Matrix Operations. In *SIGMOD*, 2007.

[9] B. Bhattacharjee et al. Efficient Index Compression in DR2. *LIW. PVLDB*, 2(2), 2009.

[10] S. Bhattacharjee et al. PShore: An Efficient Storage Framework for Managing Scientific Data. In *SSTD*, 2014.

[11] C. Huang et al. Dictionary-based Order-preserving String Compression for Main Memory Column Stores. In *SIGMOD*, 2009.

[12] M. Bilenko et al. Declarative Machine Learning - A Classification of Basic Properties and Types. *CoRR*, 2016.

[13] L. Breiman. The infinite MNIST dataset. <http://leon.breiman.org/projects/infmnist>.

[14] M. Charlier et al. Towards Estimation Error Guarantees for Distinct Values. In *SIGMOD*, 2000.

[15] R. Chittka et al. Approximate Kernel k-means: Solution to Large Scale Kernel Clustering. In *KDD*, 2011.

[16] J. Cohen et al. MAD Skills: New Analysis Practices for Big Data. *PVLDB*, 2(2), 2009.

[17] C. Constantinou and M. Lu. Quick Estimation of Data Compression and De-duplication for Large Storage Systems. In *CCP (Data Compression, Clust. and Proc.)*, 2011.

[18] G. V. Cormack. Data Compression on a Database System. *Commun. ACM*, 28(12), 1985.

[19] S. Das et al. Riscv: Integrating RISC and Hadoop. In *SIGMOD*, 2010.

[20] J. Dean and S. Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. In *OSDI*, 2004.

[21] A. Ghosh et al. SystemML: Declarative Machine Learning on MapReduce. In *ICDE*, 2011.

[22] I. J. Good. The Population Frequency of Species and the Estimation of Population Parameters. *Biometrika*, 1953.

[23] G. Graefe and L. D. Shapiro. Data Compression and Database Performance. In *Applied Computing*, 1991.

[24] P. J. Haas and L. S. Stokes. Estimating the Number of Clones in a Finite Population. *J. Amer. Statist. Assoc.*, 98(444), 1998.

[25] D. Harsh et al. Estimation of Deduplication Ratios in Large Data Sets. In *MSST (Mass Storage Sys. Tech.)*, 2012.

[26] D. Harsh et al. To Zip or not to Zip: Efficient Resource Usage for Real-Time Compression. In *PAST*, 2013.

[27] B. Huang et al. Cumbios: Optimizing Statistical Data Analysis in the Cloud. In *SIGMOD*, 2013.

[28] B. Huang et al. Resource Elasticity for Large-Scale Machine Learning. In *SIGMOD*, 2015.

[29] S. Ilieva et al. Estimating the Compression Fraction of an Index using Sampling. In *ICDE*, 2010.

[30] N. L. Johnson et al. *Transformed Discrete Distributions*. Wiley, New York, 2nd edition, 1992.

[31] 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.

[32] D. Korrer et al. SLACE: Sparse Linear Algebra in a Column-Oriented In-Memory Database System. In *SSTD*, 2014.

[33] H. Kimura et al. Compression Aware Physical Database Design. *PVLDB*, 4(10), 2011.

[34] K. Kourtis et al. Optimizing Sparse Matrix-Vector Multiplication Using Index and Value Compression. In *CF (Computing Frontiers)*, 2008.

[35] H. Lang et al. Data Blocks: Hybrid OLTP and OLAP on Compressed Storage using both Vectorization and Compression. In *SIGMOD*, 2016.

[36] P. Larson et al. SQL Server Column Store Indexes. In *SIGMOD*, 2011.

[37] M. Leisner. UCI Machine Learning Repository. Higgs, Coopers, CS Group (1990). [archive.ics.uci.edu/ml/](http://archive.ics.uci.edu/ml/).

[38] P. E. O'Neill. Model 204 Architecture and Performance. In *High Performance Transaction Systems*, 1980.

[39] Oracle. *Data Warehousing Guide, 11g Release 1*, 2007.

[40] V. Raman and G. Swart. How to Write a Table Driven Compressed Relation. In *VLDB*, 2006.

[41] V. Raman et al. DR2 with REL Acceleration: So Much More than Just a Column Store. *PVLDB*, 6(11), 2013.

[42] Y. Saito. SPARKKIT: a basic tool kit for sparse matrix computations - Version 2, 1994.

[43] M. Stonebraker et al. C-Store: A Column-oriented DBMS. In *VLDB*, 2005.

[44] M. Stonebraker et al. The Architecture of SciDB. In *SSTD*, 2011.

[45] System. *IQ 15.4 System Administration Guide*, 2013.

[46] G. Vigna and F. Valletti. Estimating the Unseen: An  $\epsilon$ -Neyman Sample Estimator for Entropy and Support Size, Shows Optimal via New CLTs. In *STOC*, 2011.

[47] F. Wotawong et al. The Implementation and Performance of Compressed Databases. *SIGMOD Record*, 29(3), 2000.

[48] S. Williams et al. Optimization of Sparse Matrix-Vector Multiplication on Emerging Multiscale Platforms. In *SC (Supercomputing Conf.)*, 2007.

[49] K. Wu et al. Optimizing Bitmap Indices With Efficient Compression. *TPDS*, 31(1), 2006.

[50] L. Ye et al. Exploiting Matrix Dependency for Efficient Distributed Matrix Computation. In *SIGMOD*, 2015.

[51] M. Zaharia et al. Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing. In *NIPS*, 2012.

[52] C. Zhang et al. Materialized Optimization for Feature Selection Workloads. In *SIGMOD*, 2014.



# 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)

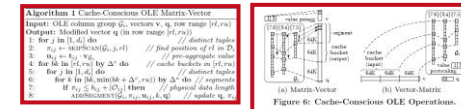
- “The rewriting for  $q=Xv$  seems wrong: To compute  $q$ , one takes each row of the matrix  $X$  and multiplies it with the vector  $v$ .”

- 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  $[r_1, r_2]$ , which guarantees disjoint results and thus avoids partial results per thread. We find  $v_j$ , the starting position of each  $v_j$  in  $D$ , via a skip scan that aggregates segment lengths until we reach  $i$  (line 5). To minimize the overhead of finding  $v_j$ , we use static scheduling (task partitioning). We further pre-compute  $u_k = v_j * v_j$  once for all tuples (line 3). For each cache-bucket of size  $\Delta$  (such that  $\Delta$  divides  $\delta$  in L2 cache, by default  $\Delta = 2\Delta$ ), we then iterate over all distinct tuples (line 5-8) but maintain the current result  $u_k$  as well. The inner non-zero  $p_i$  is written to both  $R_{ij}$  and  $R_{ji}$ . 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  $X^T \cdot X$  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  $\text{count}$ ) are efficiently computed via counts by  $\sum_{i=1}^{|G|} \sum_{j=1}^{|G|} |G_{ij}|$ . For each value, we aggregate the RLE run lengths or OLE lengths per segment, respectively. Row aggregates (e.g.,  $\text{rowsum}$ ) 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

See list at [https://pdamme.github.io/teaching/2024\\_summer/IdE/SeminarTopics.pdf](https://pdamme.github.io/teaching/2024_summer/IdE/SeminarTopics.pdf)

# Summary and Q&A



- **Course Organization, Outline, and Deliverables**
- **Structure of Scientific Papers**
- **List of Seminar Topics (Proposals)**
- **Remaining Questions?**
- **Next Lectures**
  - 02 **Scientific Reading and Writing** [Apr 29]
  - 03 **Experiments, Reproducibility, and Giving Presentations** [May 06]