

# Programmierpraktikum: Datensysteme

## 05 Background Experiments and Reproducibility

**Prof. Dr. Matthias Boehm**

Technische Universität Berlin

Berlin Institute for the Foundations of Learning and Data

Big Data Engineering (DAMS Lab)



Last update: Jun 20, 2026



# Announcements / Administrative Items



## ■ #1 Video Recording

- Hybrid lectures: in-person H 0111, zoom live streaming, video recording
- <https://tu-berlin.zoom.us/j/9529634787?pwd=R1ZsN1M3SC9BOU1OcFdmem9zT202UT09>



## ■ #2 Project Submission

- **Jul 01: Project Submissions** (DAMS instance: online submission in ISIS by 11.59pm)
- **Jul 06: Project Presentations** (in-person)

# Agenda



- **Experiments and Result Presentation**
- **Reproducibility and RDM**
- **Use Case: SliceLine Reproducibility**

# Experiments and Result Presentation

## In Computer Science (Data Management)



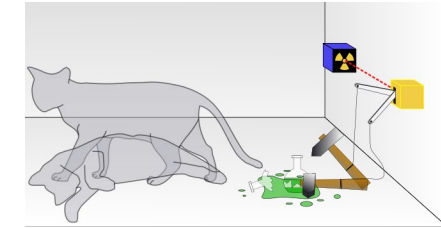
[Ioana Manolescu, Stefan Manegold:  
Performance Evaluation in Database Research:  
Principles and Experiences, **ICDE 2008**]

# Motivation



## ▪ 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 → refutable by evidence



## ▪ 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

## ▪ Good Research Fires Itself

- Initial experiments give directions for further improvements
- Problem-oriented methodology

Reproducible Research  
=  
Great, Trustworthy  
Research

# Types of Experiments



- **#1 Exploratory Experiments**
  - Tests for functional correctness
  - Unstructured experiments for initial feedback → eval feasibility
- **#2 Micro Benchmarks**
  - Measure specific aspects in controlled and understandable scope
  - Bottom-up approach
- **#3 Benchmarks**
  - Evaluate on community/own benchmarks
  - Examples: TPC-C, TPC-H, TPC-DS, JOB, MLPerf
- **#4 End-to-end Applications**
  - Evaluate in larger scope of real datasets and query workloads
  - Examples: Customer workload, ML pipelines (dataprep, training, eval)

## Programming Practical:

**#1** Micro-benchmarks for insert/get/getNext x data type performance

**#2** MIT benchmark for multi-threaded index maintenance and usage

[Additional benchmarks for skewed skews, etc]



# From Idea to Experiments

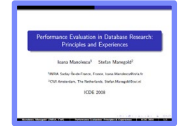
## Overview

- Proper planning helps to keep you from “getting lost”
- Repeatable experiments simplify your own work
- There is **no single way** how to **do it right**
- There are **many ways** how to **do it wrong**

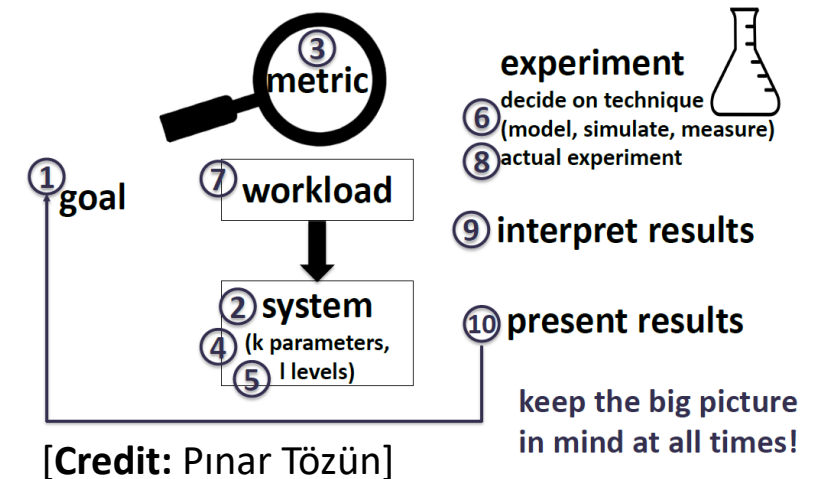
## Basic Planning

- Which data / data sets should be used?
- Which workload / queries should be run?
- Which hardware & software should be used?
- Metrics:** What to measure? How to measure?
- Comparison:** How to compare?  
**CSI:** How to find out what is going on?

[I. Manolescu, S/ Manegold: Performance Evaluation in Database Research: Principles and Experiences, **ICDE 2008**]



performance analysis



# Dataset Selection



## ■ Synthetic Data

- Generate data with specific data characteristics
- Systematic evaluation w/ datasize, sparsity, etc
- **Inappropriate for certain topics:** compression, ML accuracy

Representative  
of real data  
distributions?

## ■ “Real” Data Repositories

- Wide selection of available datasets w/ different characteristics
- UCI ML Repository: <https://archive.ics.uci.edu/ml/index.php>
- Florida Sparse Matrix Collection: <https://sparse.tamu.edu/>
- Google dataset search: <https://datasetsearch.research.google.com/>
- Common Datasets in ML: ImageNet, Mnist, CIFAR, KDD, Criteo, Adult
- Common Datasets in DM: Census, Taxi, Airlines, DBLP, benchmarks etc

Representative for  
variety of workloads /  
common case?

# Benchmarks



## ■ Overview

- Community- and organization-driven creation of agreed benchmarks
- **Benchmarks can define a field** and foster innovation

## ■ #1 Data Management

- Query processing: 007, TPC-C, TPC-E, TPC-H, TPC-DS (w/ audit)
- Join ordering: JOB

[Michael J. Carey, David J. DeWitt,  
Jeffrey F. Naughton: The oo7  
Benchmark. **SIGMOD 1993**]



[<http://www.tpc.org/tpch/>]

## ■ #2 “Big Data”

- MR/Spark: BigBench, HiBench, SparkBench
- Array Databases: GenBase

## ■ #3 Machine Learning Systems

- SLAB, DAWNbench, MLPerf, MLBench, AutoML Bench, Meta Worlds, TPCx-AI

(See AMLS course  
for details)

# Benchmarks, cont.



Closed Division Times																		
#	Submitter	System	Processor #	Accelerator #	Software	Benchmark results (minutes)							Details	Code	Notes			
						Image classification	Object detection, light-weight	Object detection, heavy-wt.	Translation, recurrent	Translation, non-recur.	Recommendation	Reinforcement Learning						
						ImageNet	COCO	COCO	WMT E-G	WMT E-G	MovieLens-20M	Go						
						ResNet-50 v1.5	SSD w/ ResNet-34	Mask-RCNN	NMT	Transformer	NCF	Mini Go						
Available in cloud																		
0.6-1	Google	TPUv3.32		TPUv3	16	TensorFlow, TPU 1.14.1.dev	42.19	12.61	107.03	12.25	10.20	[1]				details	code	none
0.6-2	Google	TPUv3.128		TPUv3	64	TensorFlow, TPU 1.14.1.dev	11.22	3.89	57.46	4.62	3.85	[1]				details	code	none
0.6-3	Google	TPUv3.256		TPUv3	128	TensorFlow, TPU 1.14.1.dev	6.86	2.76	35.60	3.53	2.81	[1]				details	code	none
0.6-4	Google	TPUv3.512		TPUv3	256	TensorFlow, TPU 1.14.1.dev	3.85	1.79		2.51	1.58	[1]				details	code	none
0.6-5	Google	TPUv3.1024		TPUv3	512	TensorFlow, TPU 1.14.1.dev	2.27	1.34		2.11	1.05	[1]				details	code	none
0.6-6	Google	TPUv3.2048		TPUv3	1024	TensorFlow, TPU 1.14.1.dev	1.28	1.21			0.85	[1]				details	code	none
Available on-premise																		
0.6-7	Intel	32x 2S CLX 8260L	CLX 8260L	64		TensorFlow						[1]	14.43			details	code	none
0.6-8	NVIDIA	DGX-1		Tesla V100	8	MXNet, NGC19.05	115.22					[1]				details	code	none
0.6-9	NVIDIA	DGX-1		Tesla V100	8	PyTorch, NGC19.05		22.36	207.48	20.55	20.34	[1]				details	code	none
0.6-10	NVIDIA	DGX-1		Tesla V100	8	TensorFlow, NGC19.05						[1]	27.39			details	code	none
0.6-11	NVIDIA	3x DGX-1		Tesla V100	24	TensorFlow, NGC19.05						[1]	13.57			details	code	none
0.6-12	NVIDIA	24x DGX-1		Tesla V100	192	PyTorch, NGC19.05			22.03			[1]				details	code	none
0.6-13	NVIDIA	30x DGX-1		Tesla V100	240	PyTorch, NGC19.05		2.67				[1]				details	code	none
0.6-14	NVIDIA	48x DGX-1		Tesla V100	384	PyTorch, NGC19.05				1.99		[1]				details	code	none
0.6-15	NVIDIA	60x DGX-1		Tesla V100	480	PyTorch, NGC19.05					2.05	[1]				details	code	none
0.6-16	NVIDIA	130x DGX-1		Tesla V100	1040	MXNet, NGC19.05	1.69					[1]				details	code	none
0.6-17	NVIDIA	DGX-2		Tesla V100	16	MXNet, NGC19.05	57.87					[1]				details	code	none
0.6-18	NVIDIA	DGX-2		Tesla V100	16	PyTorch, NGC19.05		12.21	101.00	10.94	11.04	[1]				details	code	none
0.6-19	NVIDIA	DGX-2H		Tesla V100	16	MXNet, NGC19.05	52.74					[1]				details	code	none
0.6-20	NVIDIA	DGX-2H		Tesla V100	16	PyTorch, NGC19.05		11.41	95.20	9.87	9.80	[1]				details	code	none
0.6-21	NVIDIA	4x DGX-2H		Tesla V100	64	PyTorch, NGC19.05		4.78	32.72			[1]				details	code	none
0.6-22	NVIDIA	10x DGX-2H		Tesla V100	160	PyTorch, NGC19.05						[1]				details	code	none
0.6-23	NVIDIA	12x DGX-2H		Tesla V100	192	PyTorch, NGC19.05			18.47			[1]				details	code	none
0.6-24	NVIDIA	15x DGX-2H		Tesla V100	240	PyTorch, NGC19.05		2.56				[1]				details	code	none
0.6-25	NVIDIA	16x DGX-2H		Tesla V100	256	PyTorch, NGC19.05				2.12		[1]				details	code	none
0.6-26	NVIDIA	24x DGX-2H		Tesla V100	384	PyTorch, NGC19.05				1.80		[1]				details	code	none
0.6-27	NVIDIA	30x DGX-2H, 8 chips each		Tesla V100	240	PyTorch, NGC19.05		2.23				[1]				details	code	none
0.6-28	NVIDIA	30x DGX-2H		Tesla V100	480	PyTorch, NGC19.05					1.59	[1]				details	code	none
0.6-29	NVIDIA	32x DGX-2H		Tesla V100	512	MXNet, NGC19.05	2.59					[1]				details	code	none
0.6-30	NVIDIA	96x DGX-2H		Tesla V100	1536	MXNet, NGC19.05	1.33					[1]				details	code	none



MLPerf v0.6:  
<https://mlperf.org/training-results-0-6/>,  
 MLPerf v0.7:  
<https://mlperf.org/training-results-0-7>  
 ... MLPerf v2.1 (11/2022)

96 x DGX-2H = 96 \* 16  
 = 1536 V100 GPUs  
 → ~ 96 \* \$400K = \$35M – \$40M

[<https://www.forbes.com/sites/tiriasresearch/2019/06/19/nvidia-offers-a-turnkey-supercomputer-the-dgx-superpod/#693400f43ee5>]



## ■ #1 Primary Baseline

- Existing algorithm or system infrastructure
- Main comparison point, usually with same runtime operations
- **Beware:** Avoid speedup-only results (need absolute numbers for grounding)

## ■ #2 Additional Baselines

- Alternative systems w/ different runtime and compiler
- Usually, not directly comparable but important for grounding
- E.g.,: for SystemDS → R, Julia, Spark, TensorFlow, PyTorch

## ■ Problem of **Weak Baselines**

- Authors want to show improvements
- Successive improvements over state-of-the-art don't add up

[Timothy G. Armstrong, Alistair Moffat, William Webber, **Justin Zobel**: Improvements That Don't Add Up: Ad-Hoc Retrieval Results Since 1998. **CIKM 2009**]



[Maurizio Ferrari Dacrema, Paolo Cremonesi, Dietmar Jannach: Are We Really Making Much Progress? A Worrying Analysis of Recent Neural Recommendation Approaches. **RecSys 2019**]



# Presentation – Experimental Setting



## Hardware Selection

- Multiple nodes for distributed computation
- Avoid too outdated HW (irrelevance)

## Find Balanced Level of Detail

- Underspecified: “We ran all experiments on an Intel CPU”

- Over-specified:

```
cat /proc/cpuinfo  
cat /proc/meminfo
```

```
mboehm@alpha: ~  
processor       : 111  
vendor_id      : GenuineIntel  
cpu family     : 6  
model          : 85  
model name     : Intel(R) Xeon(R)  
stepping       : 7  
microcode     : 0x5002f01  
cpu MHz        : 2201.563  
cache size    : 39424 kB  
physical id    : 1  
siblings       : 56  
core id        : 30  
cpu cores      : 28  
apicid         : 125  
initial apicid : 125  
fpu            : yes  
fpu exception  : yes  
cpuid level    : 22  
wp             : yes  
flags          : fpu vme de pse t  
pat pse36 clflush dts acpi mmx fx  
scp lm constant_tsc art arch perf  
c cpuid aperfmperf pni pclmulqdq d  
g fma cx16 xtpr pdcm pcid dca sse4
```

```
mboehm@alpha: ~  
MemTotal:      792256192 kB  
MemFree:       717562860 kB  
MemAvailable:  784507576 kB  
Buffers:       480704 kB  
Cached:        69666136 kB  
SwapCached:    0 kB  
Active:        20382624 kB  
Inactive:      50053928 kB  
Active(anon):  300800 kB  
Inactive(anon): 256 kB  
Active(file):  20081824 kB  
Inactive(file): 50053672 kB  
Unevictable:   19036 kB  
Mlocked:       19036 kB  
SwapTotal:     134217724 kB  
SwapFree:      134217724 kB  
Dirty:         116 kB  
Writeback:     0 kB  
AnonPages:     309260 kB  
Mapped:        152584 kB  
Shmem:         3124 kB  
SReclaimable:  1799996 kB  
Slab:          3003104 kB  
SReclaimable:  1799996 kB
```

## Recommendation

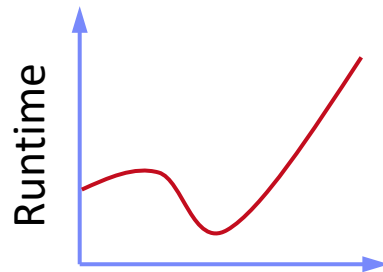
- HW components:** #nodes, CPUs, memory, network, I/O
- SW components:** OS, programming language, versions, other software
- Baselines** and configuration → Use **recent versions of baseline systems**
- Data and workloads** w/ data sizes, parameters, configurations



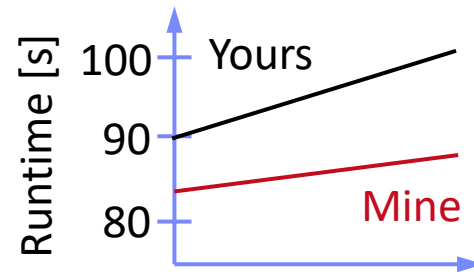
# Presentation – Figures

## ■ Axes

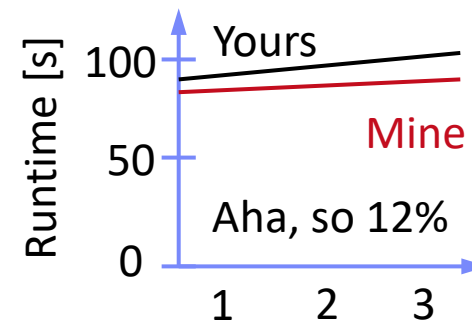
- Use **Informative axes labels with units** (e.g., Total Execution Time [ms])
- Don't cheat or mislead readers and reviewers
- **Start y-axis at 0** for linear scale (w/ log-scale impossible)



What are the units?  
Where are the ticks?



This is a misleading  
y axis

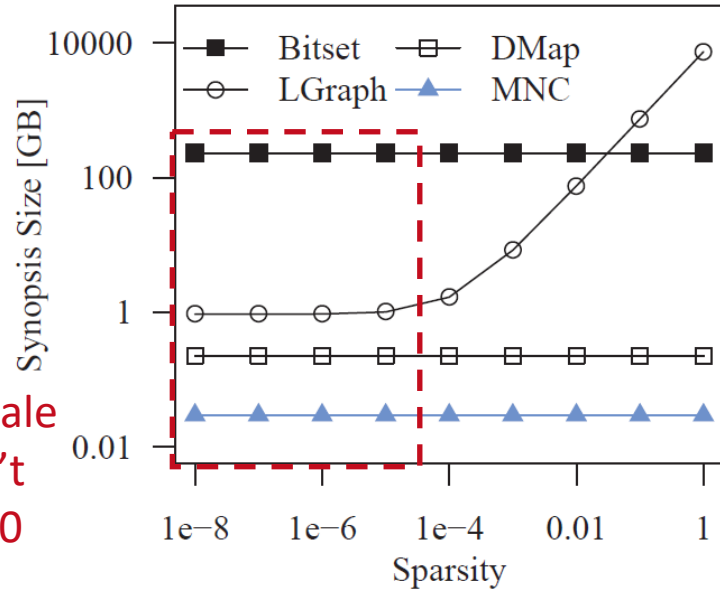


# Presentation – Figures, cont.



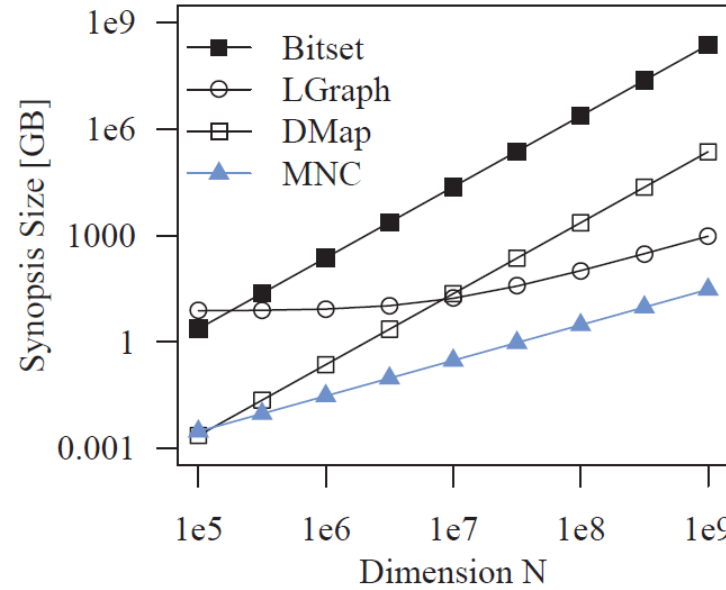
## Fair Ranges of Parameters

- Evaluate common ranges of values
- Don't hide important information → good scientific practice



For log-scale  
you can't  
start at 0

Don't limit range to  
make you look good



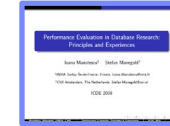
If there are multiple  
relevant parameters,  
show them all

[J. Sommer, M. Boehm, A. V. Evfimievski, B. Reinwald, P. J. Haas: MNC: Structure-Exploiting Sparsity Estimation for Matrix Expressions. **SIGMOD 2019**]



# Presentation – Figures, cont.

[I. Manolescu, S/ Manegold: Performance Evaluation in Database Research: Principles and Experiences, **ICDE 2008**]

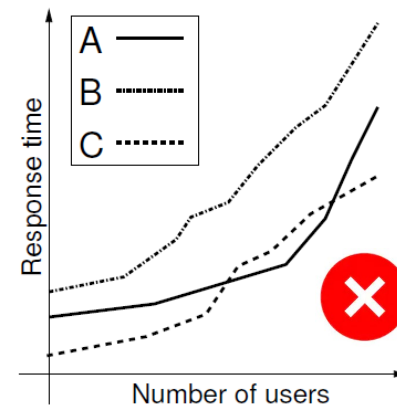
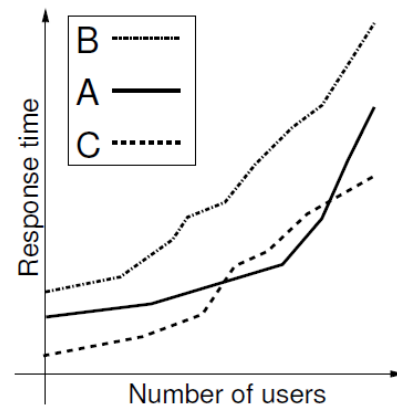
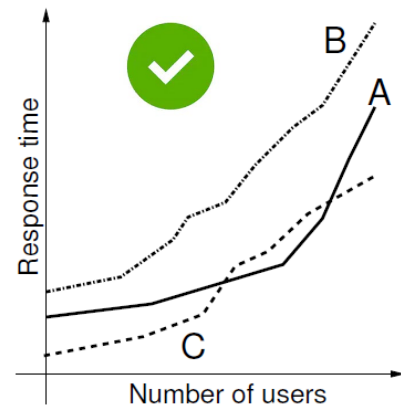
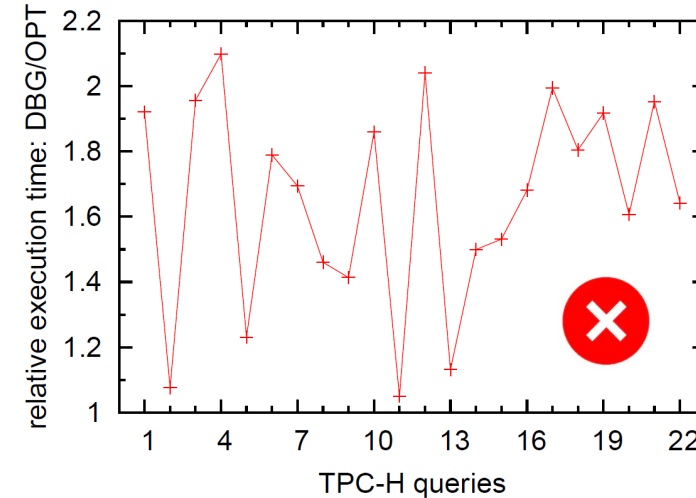


## Plots Types

- **Barplot** for categories
- **Plot + Line/linepoints** for continuous parameters
- Visible font sizes (similar to text)

## Legends

- Order them by appearance
- Attach directly to graph



Human brain is a  
**poor join processor**  
Humans **get**  
**frustrated**

# Presentation – Figures, cont.

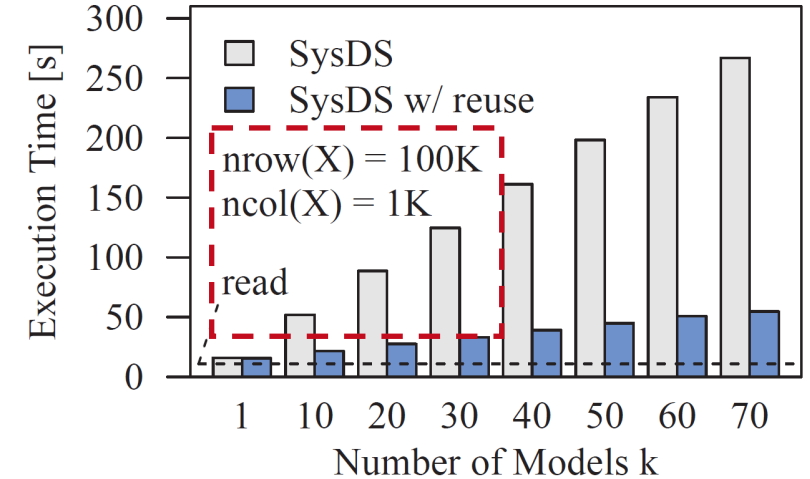
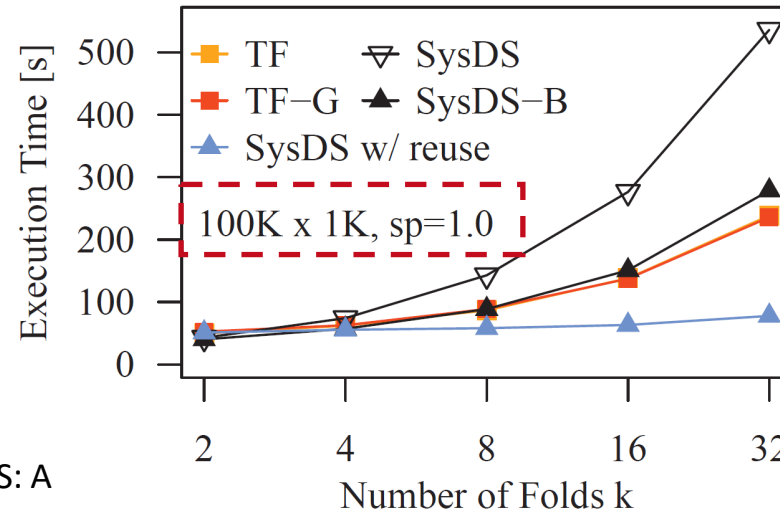


## ■ Diversity & Consistency

- **Diversity:** if applicable use mix of different plot types and tables
- **Consistency:** use consistent colors and names for same baselines

## ■ Labeling

- Make the plots self-contained
- Simplifies skimming and avoids join with text



[Matthias Boehm et al: SystemDS: A Declarative Machine Learning System for the End-to-End Data Science Lifecycle. **CIDR 2020**]

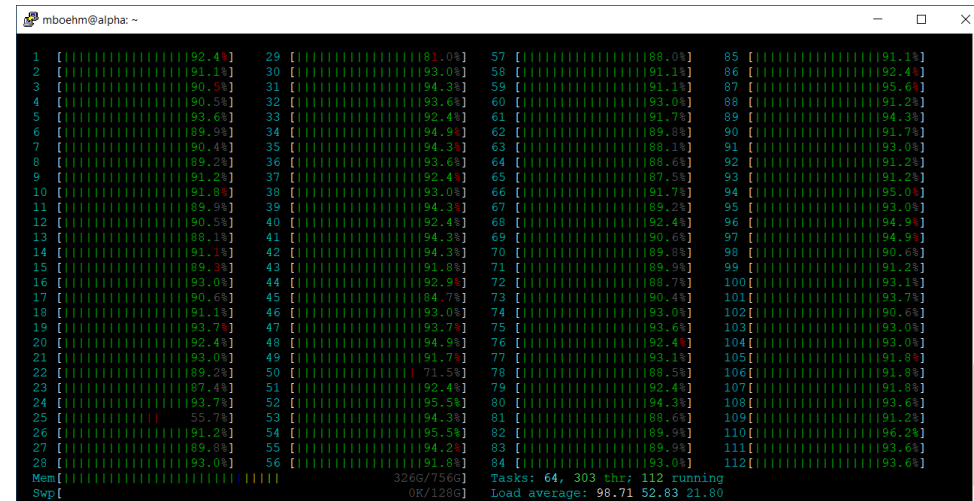


# Presentation – Result Interpretation



## Use the Right OS Tools

- System-specific tracing/statistics
- top / htop / iotop (looks **CPU bound**)
- perf -stat -d ./run.sh  
(no, it's **memory-bandwidth bound**)



Performance counter stats for './run.sh':

12721364.53 msec task-clock	#	83.640 CPUs utilized	
463352 context-switches	#	0.036 K/sec	
5455536095415 instructions	#	0.14 insn per cycle	(62.50%)
335314473273 branches	#	26.358 M/sec	(62.50%)
1463380955 branch-misses	#	0.44% of all branches	(62.50%)
2185062643097 L1-dcache-loads	#	171.763 M/sec	(62.50%)
142845949268 L1-dcache-load-misses	#	6.54% of all L1-dcache hits	(62.50%)
3375555316 LLC-loads	#	0.265 M/sec	(50.00%)
1016330404 LLC-load-misses	#	<b>30.11% of all LL-cache hits</b>	(50.00%)
152.096000108 seconds time elapsed			
12052.466691000 seconds user			
674.704421000 seconds sys			

**Don't just report the results but try to understand and explain them**

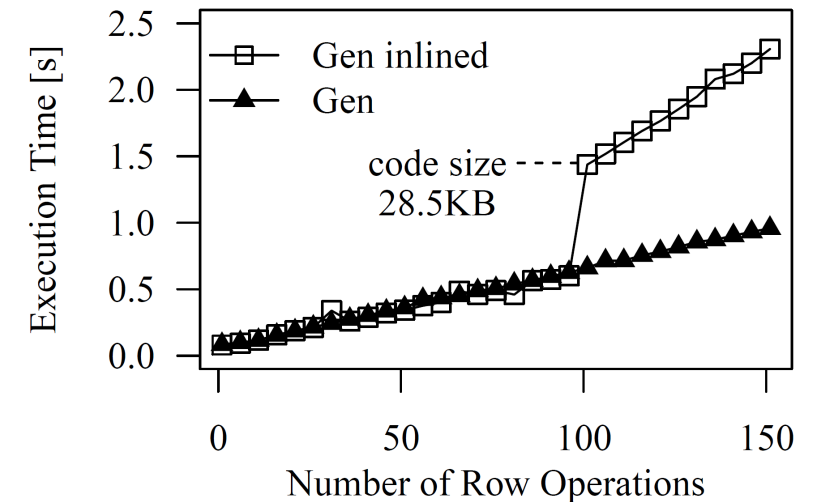
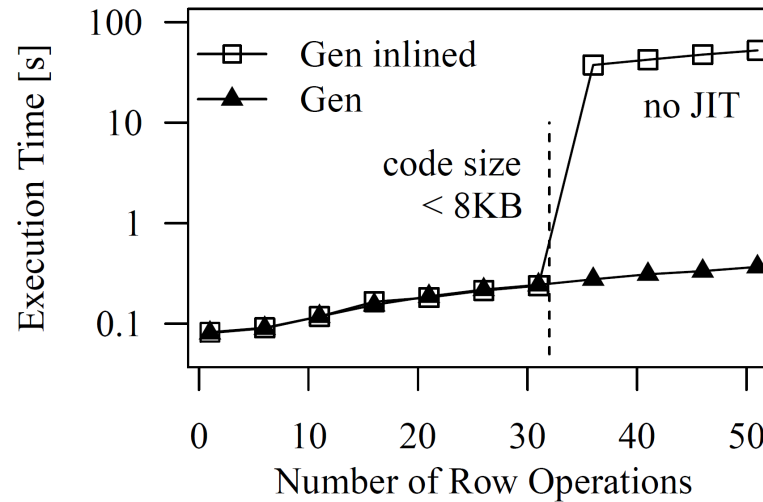
# Presentation – Result Interpretation, cont.



- Use the Right PL Tools / Flags

- E.g., Understanding Java JIT compilation  
-XX:+PrintCompilation

- E.g., Understanding HW Cache Hierarchy (L1i 32KB)  
-XX:-DontCompileHugeMethods



[Matthias Boehm et al: On Optimizing Operator Fusion Plans for Large-Scale Machine Learning in SystemML. **PVLDB 11(12) 2018**]



# Reproducibility and RDM

In Computer Science (Data Management)

# Research Data Management (RDM)



## ■ Overview

- Ensure reproducibility of research results and conclusions
- **Common problem:**
- **Create value for others** (compare, reuse, understand, extend)
- EU Projects: Mandatory proposal section & deliverable on RDM plan

**“All code and data was on the student’s laptop and the student left / the laptop crashed.”**

## ■ RDM @ TU Graz

- TU Graz RDM Policy since 12/2019, as well as faculty-specific RDM policies
- <https://www.tugraz.at/sites/rdm/home/>

“Ensure that research data, code and any other materials needed to reproduce research findings are appropriately documented, stored and shared in a research data repository in **accordance with the FAIR principles** (Findable, Accessible, Interoperable and Reusable) **for at least 10 years from the end of the research project**, unless there are valid reasons not to do so. [...] Develop a **written data management strategy** for managing research outputs within the first 12 months of the PhD study [...].”

## ■ RDM @ TU Berlin

- TU Berlin RDM Policy since 10/2019
- <https://www.tu.berlin/en/ub/szf/information-tips/what-is-research-data-management>
- [https://www.static.tu.berlin/fileadmin/www/10000000/Arbeiten/Wichtige\\_Dokumente/RDM-Policy\\_TUBerlin\\_2023\\_en.pdf](https://www.static.tu.berlin/fileadmin/www/10000000/Arbeiten/Wichtige_Dokumente/RDM-Policy_TUBerlin_2023_en.pdf)

“The **minimum storage period for research data is ten years** after either the assignment of a persistent identifier or the publication of the related work following research project completion, whichever is later.”

# FAIR Data Principles



[<https://www.go-fair.org/fair-principles/>]



## ■ #1 Findable

- Metadata and data have globally unique **persistent identifiers**
- Data describes w/ rich **meta data**; registered/indexes and searchable

## ■ #2 Accessible

- Metadata and data retrievable via open, free and universal **communication protocols**
- Metadata accessible even when data no longer available

## ■ #3 Interoperable

- Metadata and data use a formal, **accessible, and broadly applicable format**
- Metadata and data use FAIR vocabularies and qualified references

## ■ #4 Reusable

- Metadata and data described with plurality of accurate and relevant attributes
- Clear license, **associated with provenance**, meets community standards



# RDM in Practice – Example So2Sat LCZ42

[Xiao Xiang Zhu et al: So2Sat LCZ42: A Benchmark Dataset for the Classification of Global Local Climate Zones. **GRSM 2020**]

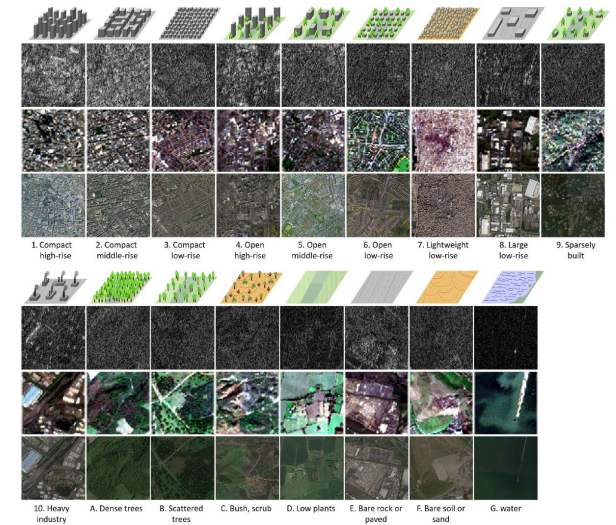
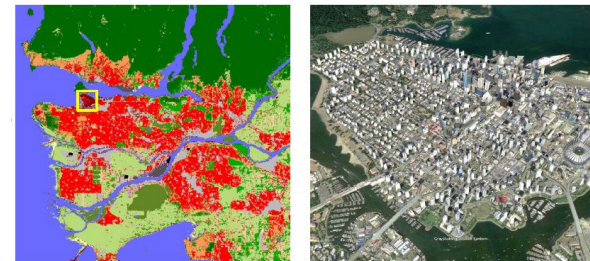


## Data and ML Pipelines

- ESA Sentinel-1/2 datasets → 4PB/year
- Training of local climate zone classifiers on So2Sat LCZ42 (15 experts, 400K instances, 10 labels each, 85% confidence, ~55GB H5)
- ML pipeline:** preprocessing, ResNet18, climate models

[So2Sat LC42 Dataset

<https://mediatum.ub.tum.de/1454690>]



## Label Creation/ Validation

- Team learning
- Labeling w/ checks
- Label validation
- Quantitative validation w/ 10 expert votes on correctness

Original Label	1	2	3	4	5	6	7	8	9	10	A	B	C	D	E	F	G
1	80	5	0	15	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	94	3	1	2	0	0	0	0	0	0	0	0	0	0	0	0
3	0	22	66	0	3	10	0	0	0	0	0	0	0	0	0	0	0
4	4	3	0	72	20	0	0	0	0	0	0	0	0	0	0	0	0
5	0	9	1	2	79	8	0	0	0	0	0	0	0	0	0	0	0
6	0	1	7	0	4	84	0	0	0	0	0	0	0	0	0	0	0
7	0	6	45	0	0	0	45	4	0	0	0	0	0	0	0	0	0
8	0	1	0	0	0	0	97	0	1	0	0	0	0	0	0	0	0
9	0	0	0	0	19	1	1	79	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	11	0	87	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	0	0	98	1	0	0	0	0	0	0	0
B	0	0	0	0	0	0	0	0	5	88	7	0	0	0	0	0	0
C	0	0	0	0	0	0	0	0	0	18	62	8	7	6	0	0	0
D	0	0	0	0	0	0	0	0	0	2	96	0	1	0	0	0	0
E	0	0	0	0	0	0	0	0	0	1	0	94	3	0	0	0	0
F	0	0	0	0	0	0	0	0	0	1	1	12	83	0	0	0	0
G	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100

Final Label (after majority voting)	1	2	3	4	5	6	7	8	9	10	A	B	C	D	E	F	G
1	80	5	0	15	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	89	4	1	4	1	1	0	0	0	0	0	0	0	0	0	0
3	0	5	70	0	1	5	18	1	0	0	0	0	0	0	0	0	0
4	5	3	0	76	15	0	0	0	0	0	0	0	0	0	0	0	0
5	0	7	1	3	85	4	0	0	0	0	0	0	0	0	0	0	0
6	0	1	6	0	5	86	0	2	0	0	0	0	0	0	0	0	0
7	0	7	32	0	0	58	2	0	0	0	0	0	0	0	0	0	0
8	0	1	0	0	0	1	95	0	2	0	0	0	0	0	0	0	0
9	0	0	0	0	0	14	1	184	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	6	0	93	0	0	0	0	0	0	0	0	0
A	0	0	0	0	0	0	0	0	97	2	0	0	0	0	0	0	0
B	0	0	0	0	0	0	0	0	3	89	7	0	0	0	0	0	0
C	0	0	0	0	0	0	0	0	0	17	66	3	7	6	0	0	0
D	0	0	0	0	0	0	0	0	0	1	3	95	0	1	0	0	0
E	0	0	0	0	0	0	0	0	0	1	0	0	95	1	0	0	0
F	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
G	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100



# RDM in Practice – Example So2Sat LCZ42, cont.



**[So2Sat LCZ42 Dataset]**  
<https://mediatum.ub.tum.de/1454690>

mediaTUM Gesamtbestand » Forschungsdaten » Datenverarbeitung, Informatik »

Document type: Forschungsdaten  
Publication date: 28.09.2018  
Authors: Zhu, Xiaoxiang ; Hu, Jingliang ; Qiu, Chunping ; Shi, Yilei ; Bagheri, Hossein ; Kang, Jian ; Li, Hao ; Mou, Lichao ; Zhang, Guicheng ; Häberle, Matthias ; Han, Shiyao ; Hua, Yuansheng ; Huang, Rong ; Hughes, Lloyd ; Sun, Yao ; Schmitt, Michael ; Wang, Yuanyuan  
Author affiliation: Zhu, Xiaoxiang (TUM, DLR); Hu, Jingliang (DLR); Qiu, Chunping (TUM); Shi, Yilei (TUM); Bagheri, Hossein (TUM); Kang, Jian (TUM); Li, Hao (TUM); Mou, Lichao (TUM); Zhang, Guicheng (TUM); Häberle, Matthias (DLR); Han, Shiyao (TUM); Hua, Yuansheng (TUM); Huang, Rong (TUM); Hughes, Lloyd (TUM); Sun, Yao (DLR); Schmitt, Michael (TUM); Wang, Yuanyuan (TUM)  
Publisher: TUM  
Title: So2Sat LCZ42  
Identifier: doi:10.14459/2018MP1454690  
Time of production: 30.08.2018  
Subject area: BAU Bauingenieurwesen, Vermessungswesen; DAT Datenverarbeitung, Informatik; GEO Geowissenschaften  
Resource type: Abbildungen von Objekten / image of objects  
Data type: Bilder / images  
Description: So2Sat LCZ42 is a dataset consisting of co-registered synthetic aperture radar and multispectral optical image patches acquired by the Sentinel-1 and Sentinel-2 remote sensing satellites, and the corresponding local climate zones (LCZ) label. The dataset is distributed over 42 cities across different continents and cultural regions of the world.  
Method of data assessment: Sentinel-1 image downloaded from ESA SciHub, and prepared by ESA SNAP software. Sentinel-2 image semi-automatically downloaded and prepared using Google Earth Engine and MATLAB. The local climate zones labels were manually labeled.  
Links: This is version 1 of the dataset. See download link below  
Link to version 2 (containing additional test data): <https://mediatum.ub.tum.de/1483140>  
Link to version 3 (containing 2 more training testing splits): <https://mediatum.ub.tum.de/1613658>  
A detailed description of this dataset, and its different versions: <https://github.com/zhu-xlab/So2Sat-LCZ42>  
Key words: local climate zones ; big data ; classification ; remote sensing ; deep learning ; data fusion ; synthetic aperture radar imagery ; optical imagery  
Technical remarks: **View and download (51.8 GB, 6 files)**  
The data server also offers downloads with FTP  
The data server also offers downloads with rsync (password m1454690):  
rsync rsync://m1454690@dataserv.ub.tum.de/m1454690/  
Language: en  
Rights: by, <http://creativecommons.org/licenses/by/4.0>  
Horizon 2020: ERC-2016-StG-714087

Key words: local climate zones ; big data ; classification ; remote sensing ; deep learning ; data fusion ; synthetic aperture radar imagery ; optical imagery

Technical remarks: **View and download (51.8 GB, 6 files)**  
The data server also offers downloads with FTP  
The data server also offers downloads with rsync (password m1454690):  
rsync rsync://m1454690@dataserv.ub.tum.de/m1454690/  
Language: en  
Rights: by, <http://creativecommons.org/licenses/by/4.0>  
Horizon 2020: ERC-2016-StG-714087

Occurrences:

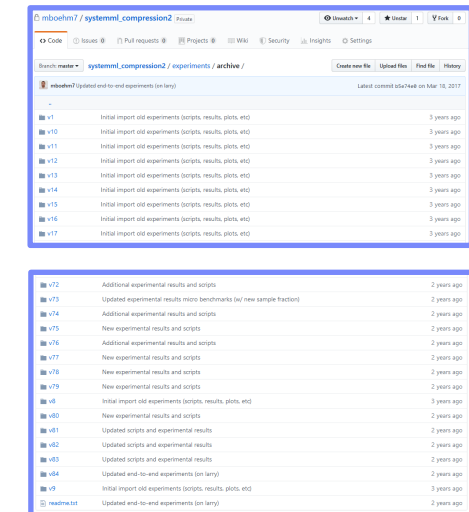
Name	Size	Modified
chvebsamechul02	... + 1.00	3 years ago
LCZ42	... 19.00	3 years ago
read_file.py	... 1.00	3 years ago
Readme.txt	... 3.00	3 years ago
README.txt	... 3.00	3 years ago
training10	... 48.4 GB	3 years ago
validation10	... 3.3 GB	3 years ago

## ■ Code and Artifacts

- Apache SystemDS: <https://github.com/apache/systemds> (OSS)
  - Complete code history, src/bin releases (SystemDS 3.1.0 in Mar/2023)
  - DIA / AMLS programming projects in SystemDS
- Additional private github repos for student projects / prototypes

## ■ Central Paper Repository

- All paper submissions w/ latex sources, figures, reviews, rebuttals, etc
  - All paper-related experiments
    - Archive: append-only experimental results
    - Plots: scripts and figures of plots
    - Results: latest results used for the current plots
    - Scripts: data preparation, baselines, benchmarks
- ➔ Automate your experiments as much as possible



# SIGMOD Reproducibility Process



## ■ Overview

- Accepted papers can submit package, verified by committee
- ACM Results Replicated / ACM Artifacts Available labels
- Most Reproducible Paper Award (\$750, visibility)

## ■ #1 Replicability (aka **Repeatability**)

- Recreate result data and graphs shown in the final paper
- **Expected:** same trend of baseline comparisons, parameter influence



## ■ #2 **Reproducibility**

- Verify robustness of results wrt parameters and environments
- **Examples:** different data and workload characteristics, hardware



## ■ Ideal Reproducibility Submission

“At a minimum the authors should provide a complete set of scripts to install the system, produce the data, run experiments and produce the resulting graphs along with a detailed Readme file that describes the process step by step so it can be easily reproduced by a reviewer. **The ideal reproducibility submission consists of**

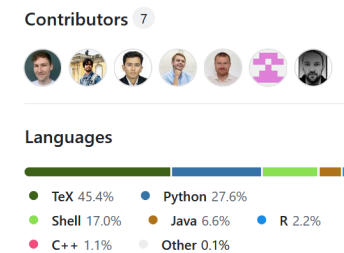
1. installs all systems needed,
2. generates or fetches all needed input data,
3. reruns all experiments and generates all results,
4. generates all graphs and plots, and finally,
5. recompiles the sources of the paper

... to produce a new PDF for the paper that contains the new graphs. “

[Credit:  
[db-reproducibility.seas.harvard.edu/  
#Guidelines](https://db-reproducibility.seas.harvard.edu/#Guidelines)]

## ■ DAMS Reproducibility Repository

- Full reproducibility packages of (almost) all research papers
- <https://github.com/damslab/reproducibility>



## Excursus: SIGMOD Contributions Award 2020



The **SIGMOD 2020 Contributions Award** recognizes the innovative work in the data management community to **encourage scientific reproducibility** of our publications. Reproducibility was introduced at the 2008 SIGMOD Conference and since then has influenced how the community approaches experimental evaluation.

**Philippe Bonnet**  
(ITU Copenhagen)

**Stratos Idreos**  
(Harvard)

**Dennis Shasha**  
(NYU)



**Ioana Manolescu**  
(Ecole Polytechnique)

**Stefan Manegold**  
(CWI Amsterdam)

**Juliana Freire**  
(NYU)

# Tools for Automation



## ■ Motivation

- Tooling for running artifacts in a self-contained manner
- Metadata storage in semi-structured formats like JSON/XML

## ■ Examples

Name	Target	Link
CK	ML Systems	<a href="http://ctuning.org">http://ctuning.org</a>
CWL	Analysis Workflows	<a href="http://commonwl.org">http://commonwl.org</a>
Popper	Container Workflows	<a href="https://github.com/systemslab/popper">https://github.com/systemslab/popper</a>
ReproZip	General-purpose Bundles	<a href="http://reprozip.org">http://reprozip.org</a>
Sciunit	Self-contained Experiments Bundles	<a href="http://sciunit.run">http://sciunit.run</a>
Sumatra	Numerical Simulations	<a href="https://github.com/open-research/sumatra">https://github.com/open-research/sumatra</a>

# Use Case: SliceLine Reproducibility



[Svetlana Sagadeeva, Matthias Boehm:  
SliceLine: Fast, Linear-Algebra-based Slice Finding for ML Model Debugging, **SIGMOD 2021**]

[<https://github.com/damslab/reproducibility/tree/master/sigmod2021-sliceline-p218>]



sliceline

[Credit: sliceline,  
Silicon Valley, HBO]

# Overview SliceLine Reproducibility Scripts



## ▪ ./runAll.sh

```
#!/bin/bash

# clean original results
rm -R results/*;
rm -R plots/*;

# setup, run experiments, plots
./run1SetupDependencies.sh;
./run2SetupSystemDS.sh;
./run3DownloadData.sh;
./run4PrepareLocalData.sh; # on scale-up node
./run5LocalExperiments.sh; # on scale-up node
./run6PrepareDistData.sh; # on scale-out cluster
./run7DistExperiments.sh; # on scale-out cluster
./run8PlotResults.sh;
```

# 1) Setup Dependencies and Apache SystemDS



```
./run1SetupDependencies.sh;  
./run2SetupSystemDS.sh;  
./run3DownloadData.sh;  
./run4PrepareLocalData.sh;  
./run5LocalExperiments.sh;  
./run6PrepareDistData.sh;  
./run7DistExperiments.sh;  
./run8PlotResults.sh;
```

## ■ Dependencies

```
sudo apt update  
sudo apt install -y openjdk-8-jdk-headless  
sudo apt install -y maven  
sudo apt install -y git  
sudo apt install -y r-base  
export JAVA_HOME=/usr/lib/jvm/java-1.8.0-openjdk-amd64  
sudo Rscript exp/setup/installDependencies.R
```

## ■ Apache SystemDS

```
# clone Apache SystemDS repository  
rm -rf systemds #cleanup  
git clone https://github.com/apache/systemds.git  
# checkout commit hash as of camera-ready version  
cd systemds  
git checkout -b reproducibility 627825c25d5a5938a772a78ce037c57e68611998  
# build systemds and prepare all dependencies  
mvn clean package -P distribution
```

## 2) Download Data

- **Download and Simple Preprocessing**

```
./run1SetupDependencies.sh;  
./run2SetupSystemDS.sh;  
./run3DownloadData.sh;  
./run4PrepareLocalData.sh;  
./run5LocalExperiments.sh;  
./run6PrepareDistData.sh;  
./run7DistExperiments.sh;  
./run8PlotResults.sh;
```



### # Adult

```
curl https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data -o data/Adult.csv;  
sed -i '$d' data/Adult.csv; # fix empty line at end of file
```

### # KDD'98

```
curl https://archive.ics.uci.edu/ml/.../epsilon_mirror/cup98lrn.zip -o data/cup98lrn.zip;  
unzip data/cup98lrn.zip -d data;  
mv data/cup98LRN.txt data/KDD98.csv;  
rm data/cup98lrn.zip;  
sed -i 's/-/ /g' data/KDD98.csv; # fix suffix - at 5th column (numerical)  
# note: workaround for macOS issue: sed -i '' 's/-/ /g' data/KDD98.csv;
```

### # CriteoD21 (this one might take a while, only needed for distributed experiments)

```
curl http://azuremlsampleexperiments.blob.core.windows.net/criteo/day_21.gz -o data/Criteo_D21.gz;  
gzip -d data/Criteo_D21.gz;  
mv data/Criteo_D21 data/Criteo_D21.csv;
```

## 5) Local Experiments

- `run5LocalExperiments.sh` → run Experiment 1, 2, 3, 4a, 4b, 4c, 5
- `runExperiment1.sh`

```
./run1SetupDependencies.sh;  
./run2SetupSystemDS.sh;  
./run3DownloadData.sh;  
./run4PrepareLocalData.sh;  
./run5LocalExperiments.sh;  
./run6PrepareDistData.sh;  
./run7DistExperiments.sh;  
./run8PlotResults.sh;
```



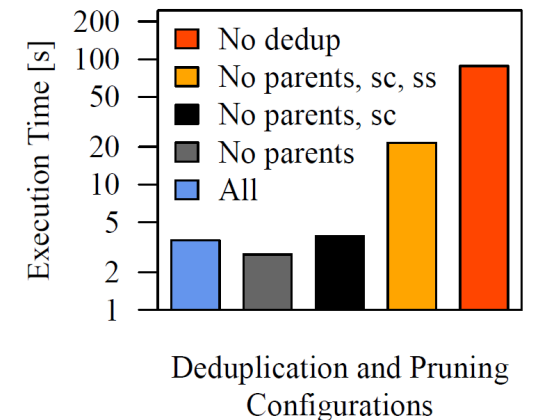
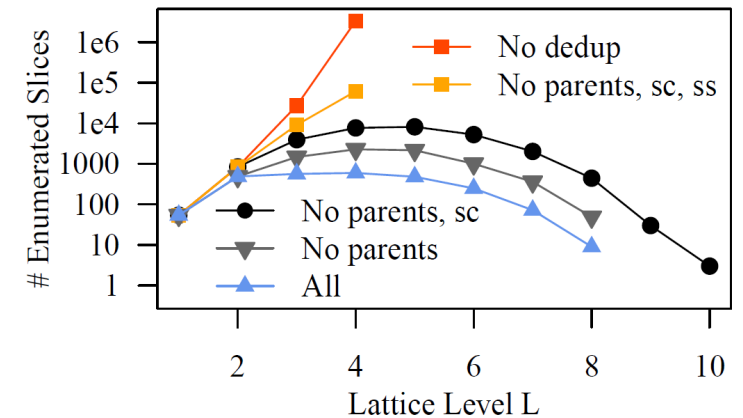
```
#!/bin/bash  
CMD="java -Xmx600g -Xms600g -cp ./lib/*:./SystemDS.jar org.apache.sysds.api.DMLScript "  
  
for config in 1 2 3 4 5  
do  
  for rep in {1..3} #for all repetitions  
  do  
    start=$(date +%s%N)  
    $CMD -f exp/explocal/SlicingExp1_${config}.dml -exec singlenode -stats \  
      -args data/Salaries_X.csv data/Salaries_o_e.csv ${config} results/Experiment1_p${config}.dat  
    end=$(date +%s%N)  
    echo ${config}", "$(((end-$start) / 1000000))" >> results/Experiment1_times.dat  
  done  
done
```

## 8) Plot Results

- **Calling R Scripts** `Rscript exp/plotting/Experiment1a.r;`  
`Rscript exp/plotting/Experiment1b.r;`

- **Plotting**  
`pdf(file="plots/Experiment1a.pdf",`  
`width=5, height=4.0, family="serif", pointsize=14)`  
  
`data1 = ...`  
`plot_colors <- c("cornflowerblue","gray40","black","orange","orangered")`  
  
`plot(points, data3,`  
`type="o", pch=19, cex=1.1, col=plot_colors[3], ylim = c(0.5,3500000),`  
`xlab="", ylab="", axes=FALSE, bg=plot_colors[3], log="y", lwd=1.1, lty=1)`  
`legend("topright",`  
`c("No dedup","No parents, sc, ss"), col=plot_colors[5:4],`  
`pch=c(15,15), lty=c(1), lwd=c(1.1), bty="n");`  
`legend( 1.5, 600,`  
`c("No parents, sc","No parents", "All"), col=plot_colors[3:1],`  
`pt.bg=plot_colors[3:1], pch=c(19,25,17), lty=c(1), lwd=c(1.1), bty="n");`  
  
`dev.off()`

```
./run1SetupDependencies.sh;  
./run2SetupSystemDS.sh;  
./run3DownloadData.sh;  
./run4PrepareLocalData.sh;  
./run5LocalExperiments.sh;  
./run6PrepareDistData.sh;  
./run7DistExperiments.sh;  
./run8PlotResults.sh;
```



# Parting Thoughts: Conflicting Questions (and Answers)



## ▪ #1 Larger System Development

- Your research does not benefit from a general improvement
- New work by others creates more work for you



**Always do the “right thing”™,  
research is a long game and  
persistently attached to your name**

## ▪ #2 What if Results are Unconvincing

- Despite effort baselines are sometimes better
- Can influence by baselines, workload, presentation
- **Several suicide cases**



**Talk to your Supervisor (and or  
trusted BIFOLD persons);  
and always stay ethical**

## ▪ #3 What if Results were Wrong

- Creating the reproducibility package after acceptance can reveal bugs in the original experiments
- Should you rerun all experiments



**Talk to your Supervisor,  
do the “right thing”™,  
yes always admit mistakes  
(camera-ready, errata, website)**

## Summary & QA



# Thanks

- Experiments and Result Presentation
- Reproducibility and RDM
- Use Case: SliceLine Reproducibility
  
- Upcoming Dates
  - Jul 01: **Project Submissions** (DAMS instance: online submission in ISIS by 11.59pm)
  - Jul 06: **Project Presentations** (in-person)