

# Comparative Recommender System Evaluation: Benchmarking Recommendation Frameworks

Alan Said

@alansaid

TU Delft

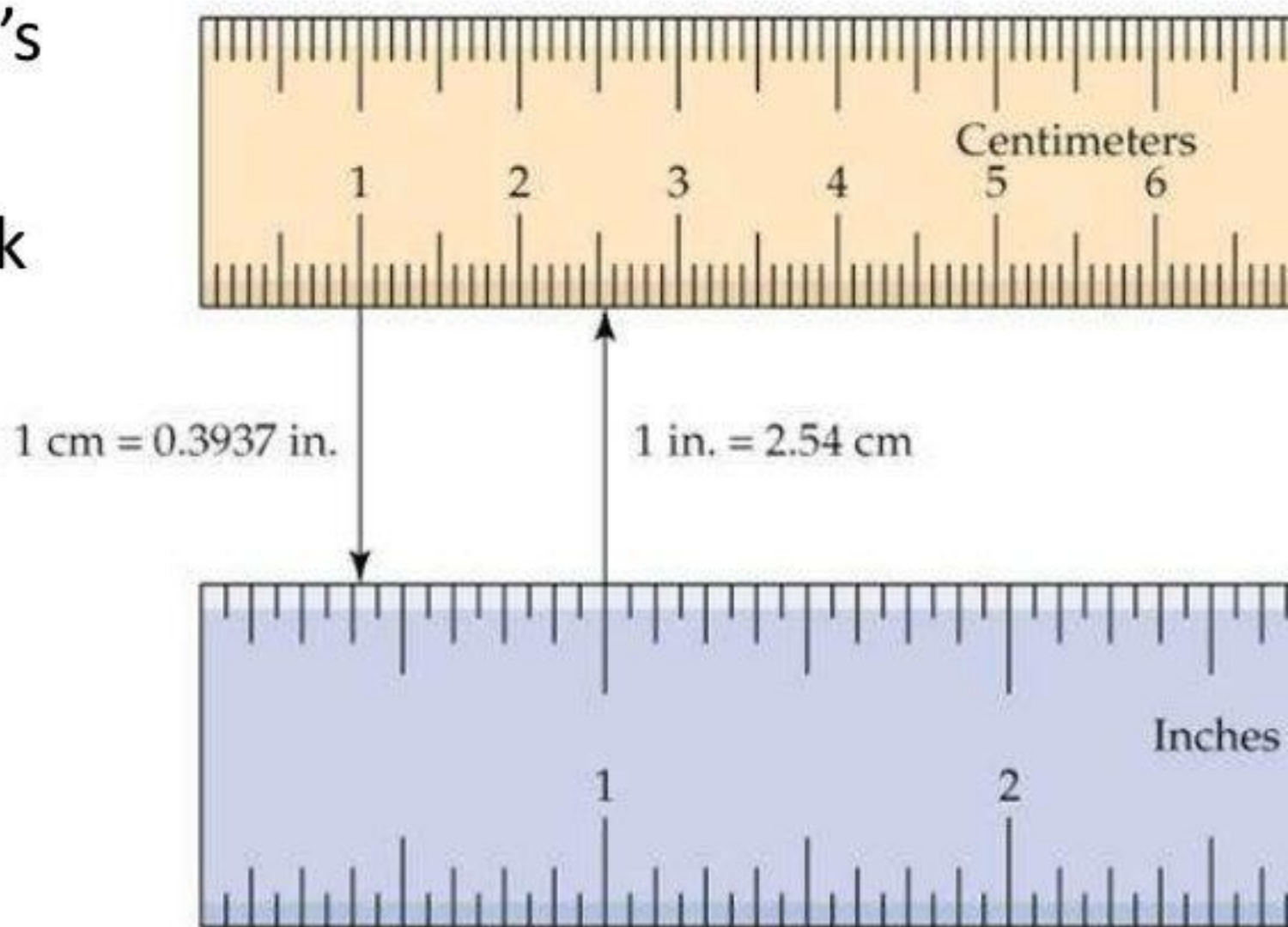
Alejandro Bellogín

@abellogin

Universidad Autónoma de Madrid

# A RecSys paper outline

- We have a new model – it's great
- We used %DATASET% 100k to evaluate it
- It's 10% better than our baseline
- It's 12% better than [Authors, 2010]







# Benchmarking





Python-recsys

Recommendable



ensKit



**mahout**



W:

mrec

LibRec



**SLIM**

**MyMedia**

**SVDfeature**





## What are the differences?

Some things just work differently

- Data splitting
- Algorithm design (implementation)
- Algorithm optimization
- Parameter values
- Evaluation
- Relevance/ranking
- Software architecture
- etc

**Different design choices!!**

# How do these choices affect evaluation results?



# Evaluate evaluation

- Comparison of frameworks
- Comparison of implementation
- Comparison of results
- Objective benchmarking



# Algorithmic Implementation

Framework	Class	Similarity	
		<b>Item-based</b>	
LensKit	ItemItemScorer	CosineVectorSimilarity, PearsonCorrelation	
Mahout	GenericItemBasedRecommender	UncenteredCosineSimilarity, PearsonCorrelationSimilarity	
MyMediaLite	ItemKNN	Cosine, Pearson	
		<b>User-based</b>	<b>Parameters</b>
LensKit	UserUserItemScorer	CosineVectorSimilarity, PearsonCorrelation	SimpleNeighborhoodFinder, NeighborhoodSize
Mahout	GenericUserBasedRecommender	UncenteredCosineSimilarity, PearsonCorrelationSimilarity	NearestNUserNeighborhood, neighborhoodsize
MyMediaLite	UserKNN	Cosine, Pearson	neighborhoodsize
		<b>Matrix Factorization</b>	
LensKit	FunkSVDItemScorer	IterationsCountStoppingCondition, factors, iterations	
Mahout	SVDRecommender	FunkSVDFactorizer, factors, iterations	
MyMediaLite	SVDPlusPlus	factors, iterations	



# There's more than algorithms though

There's the data, evaluation, and more

## Data splits

- 80-20 Cross-validation
- Random Cross-validation
- User-based cross validation
- Per-user splits
- Per-item splits
- Etc.

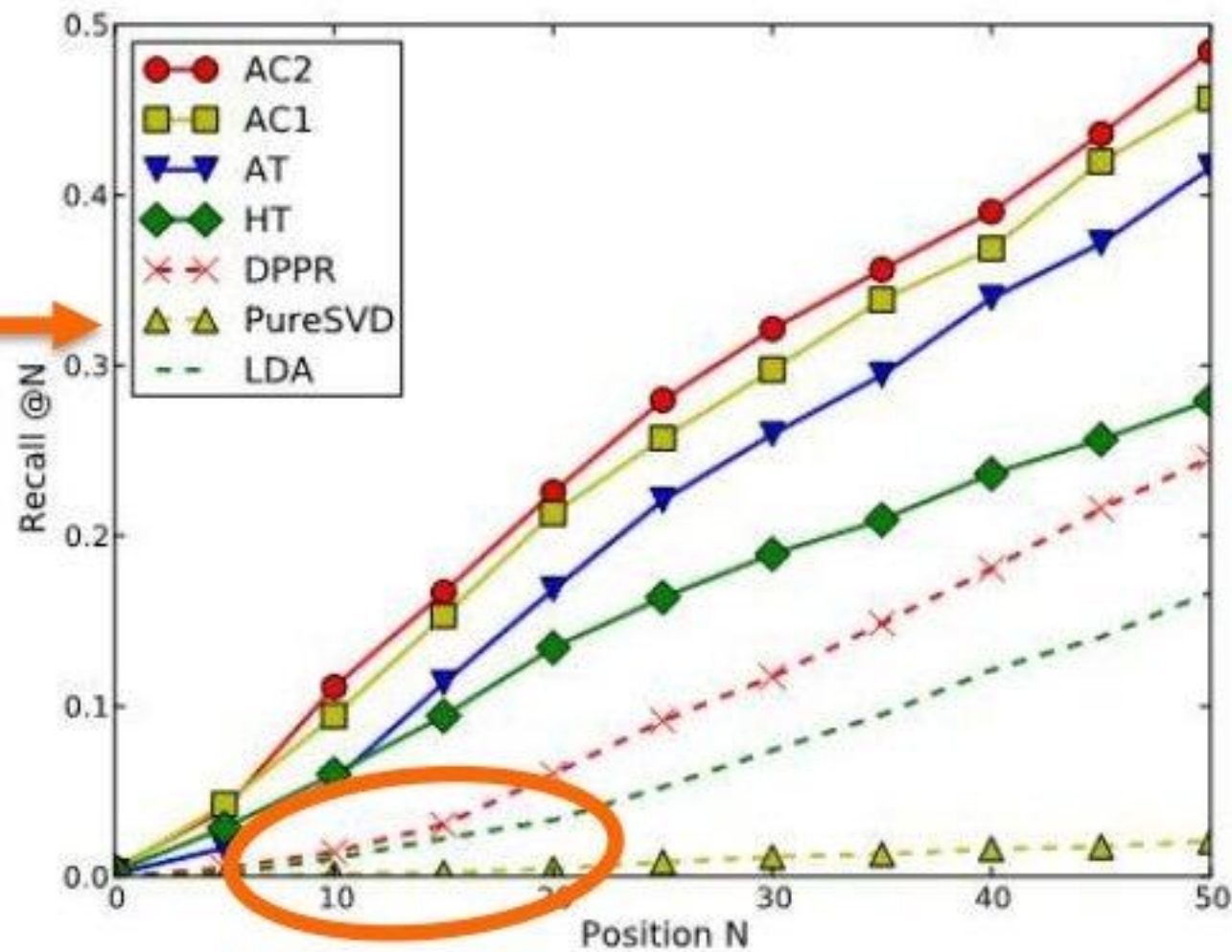
## Evaluation

- Metrics
- Relevance
- Strategies

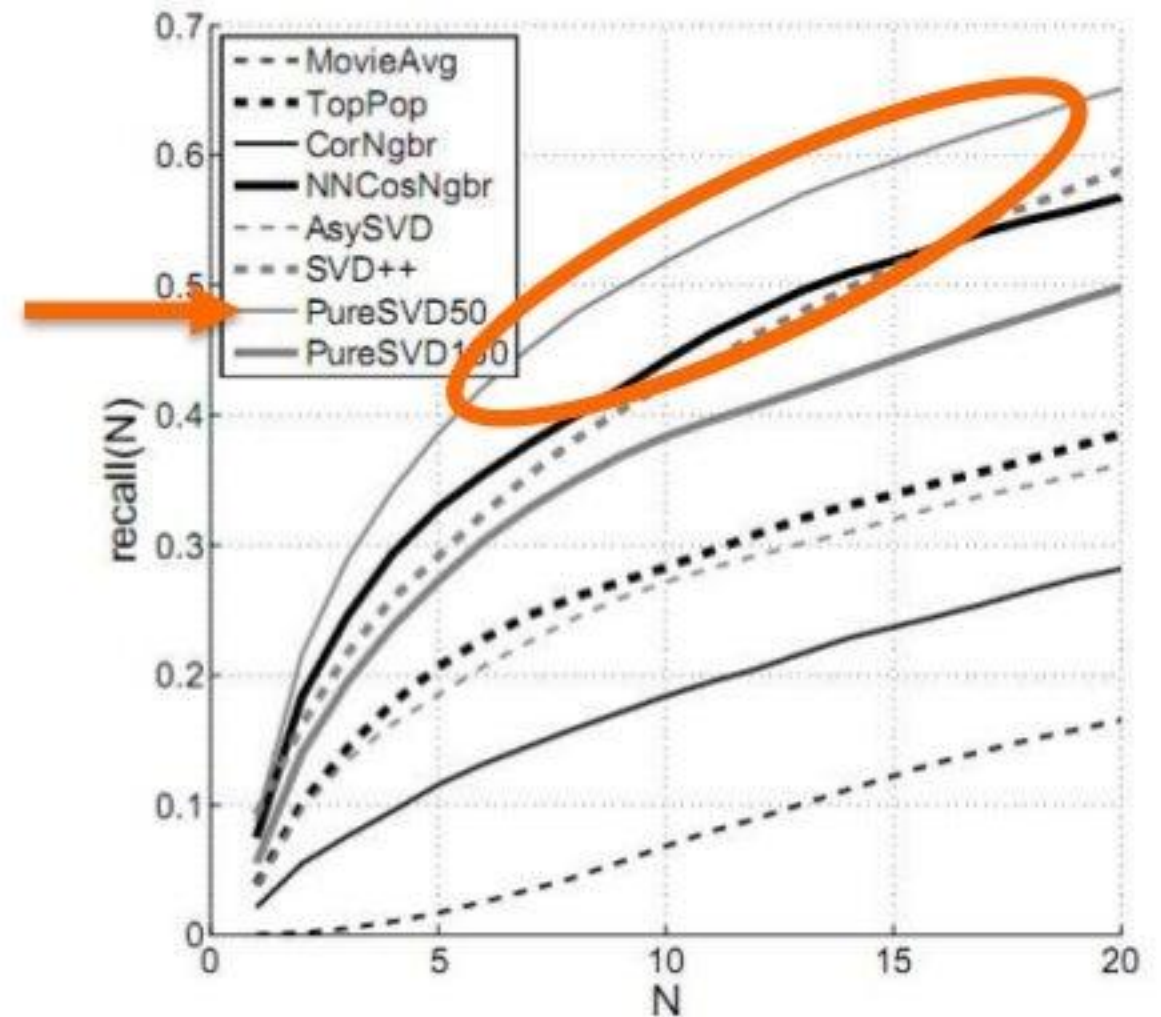




# Real world examples



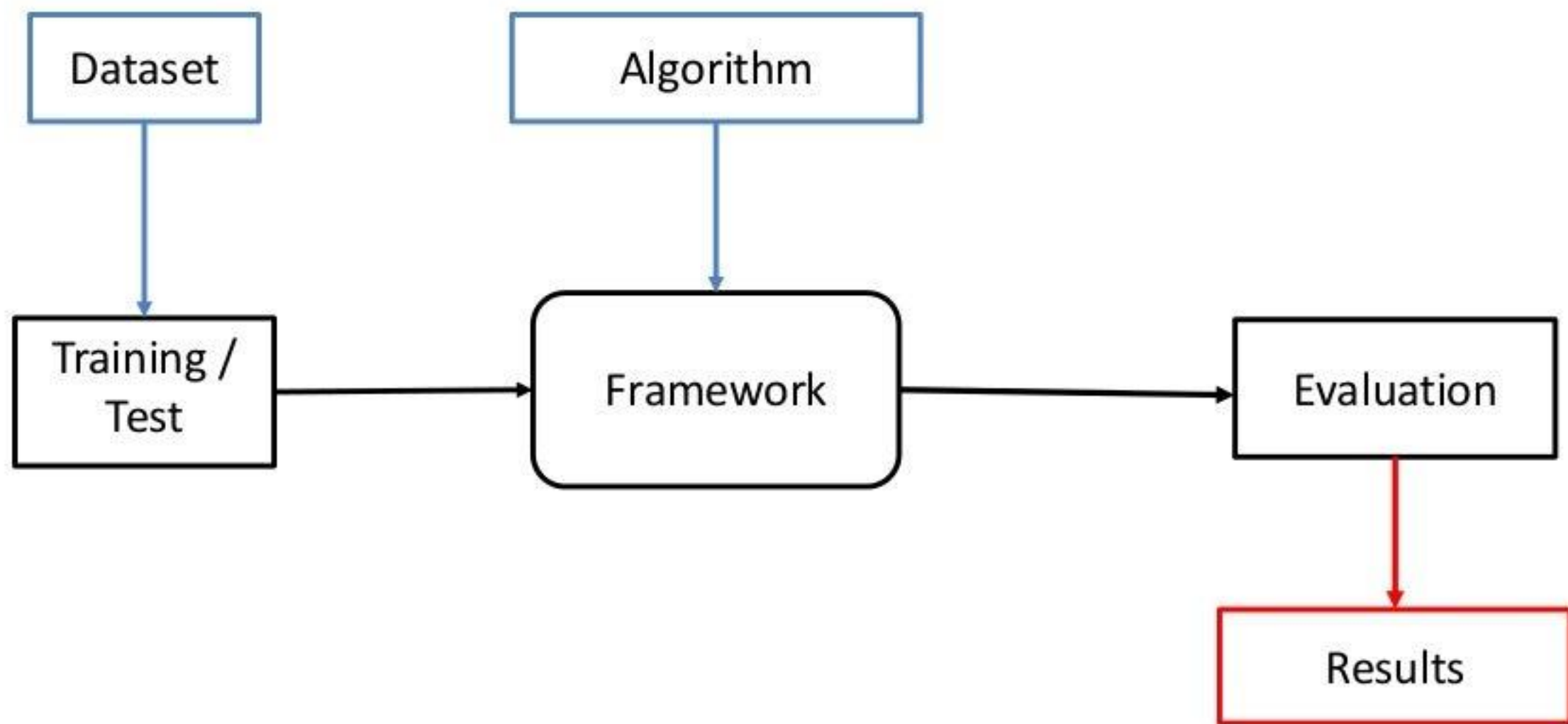
MovieLens 1M  
[Yin et al, 2012]



(a) recall  
MovieLens 1M  
[Cremonesi et al, 2010]

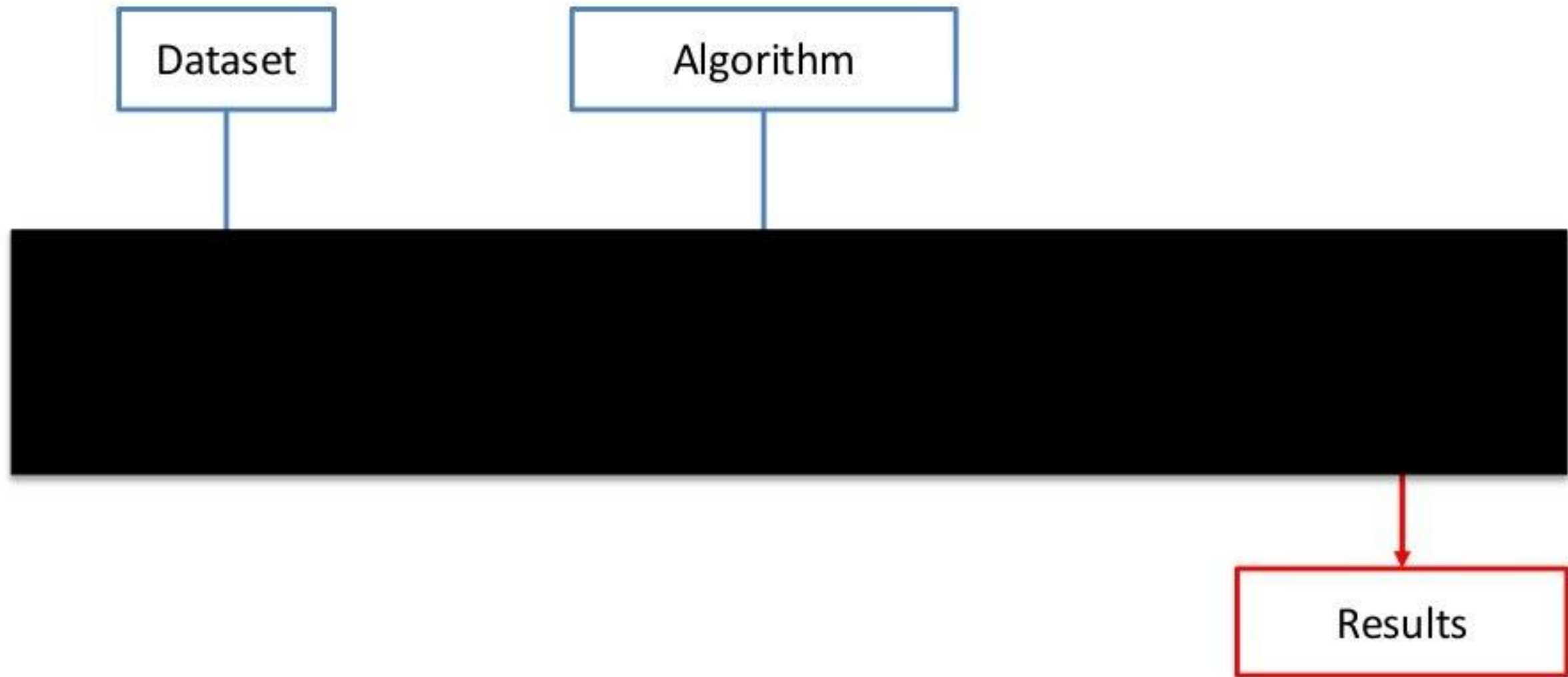


# Evaluation





# Internal Evaluation





# Internal Evaluation Results

Algorithm	Framework	nDCG
IB Cosine	Mahout	0,00041478
IB Cosine	Lenskit	0,94219205
IB Pearson	Mahout	0,00516923
IB Pearson	Lenskit	0,92454613
SVD50	Mahout	0,10542729
SVD50	Lenskit	0,94346409
UB Cosine	Mahout	0,16929545
UB Cosine	Lenskit	0,94841356
UB Pearson	Mahout	0,16929545
UB Pearson	Lenskit	0,94841356



Algorithm	Framework	RMSE
IB Cosine	Lenskit	1,01390931
IB Cosine	MyMediaLite	0,92476162
IB Pearson	Lenskit	1,05018614
IB Pearson	MyMediaLite	0,92933246
SVD50	Lenskit	1,01209290
SVD50	MyMediaLite	0,93074012
UB Cosine	Lenskit	1,02545490
UB Cosine	MyMediaLite	0,93419026

A brass balance scale is shown on a blue surface against a purple wall. The scale is slightly tilted, with the left pan lower than the right. The text "We need a fair and common evaluation protocol!" is overlaid in white.

We need a fair and common evaluation protocol!



# Reproducible evaluation - Benchmarking

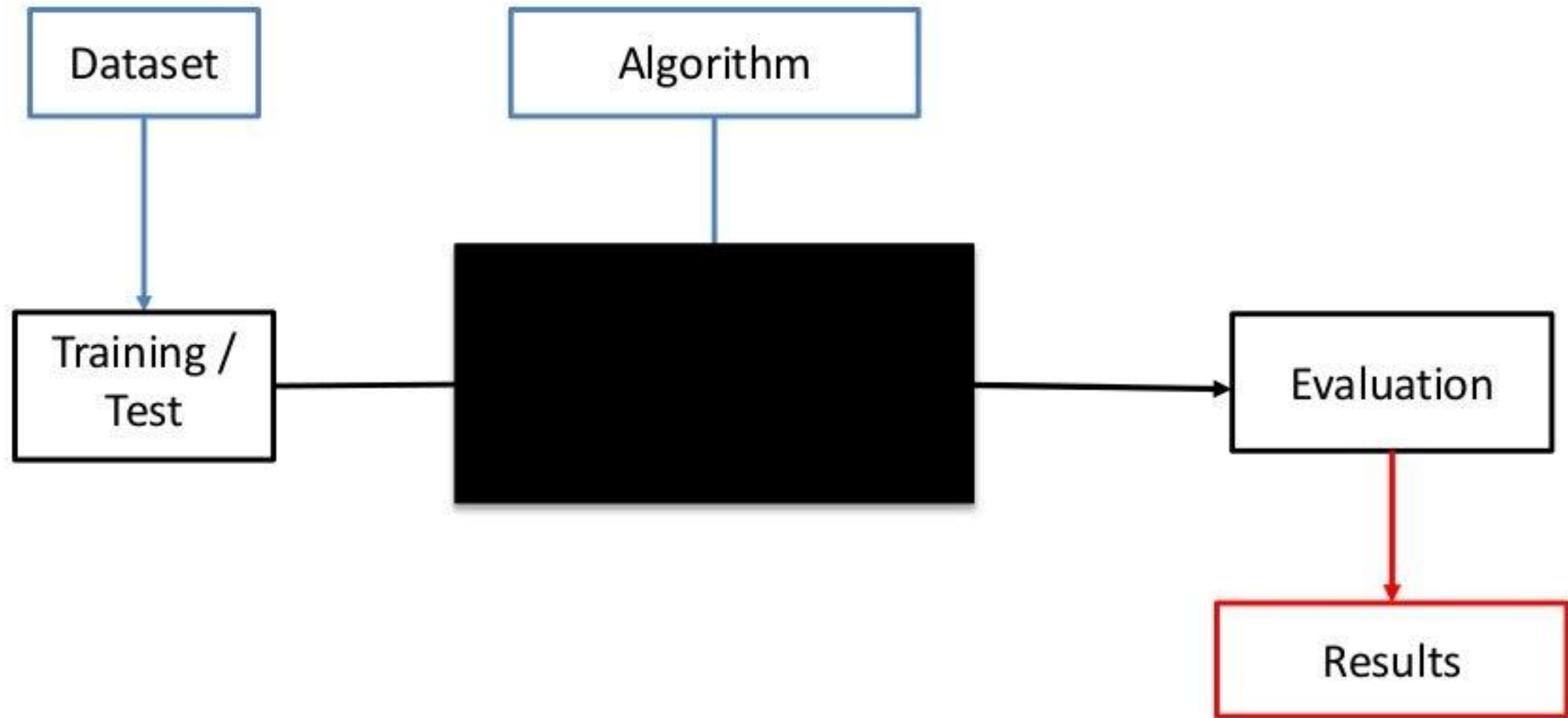
Control all parts of the process

- Data Splitting strategy
- Recommendation (black box)
- Candidate items generation (what items to test)
- Evaluation

<http://rival.recommenders.net>



# Controlled Evaluation





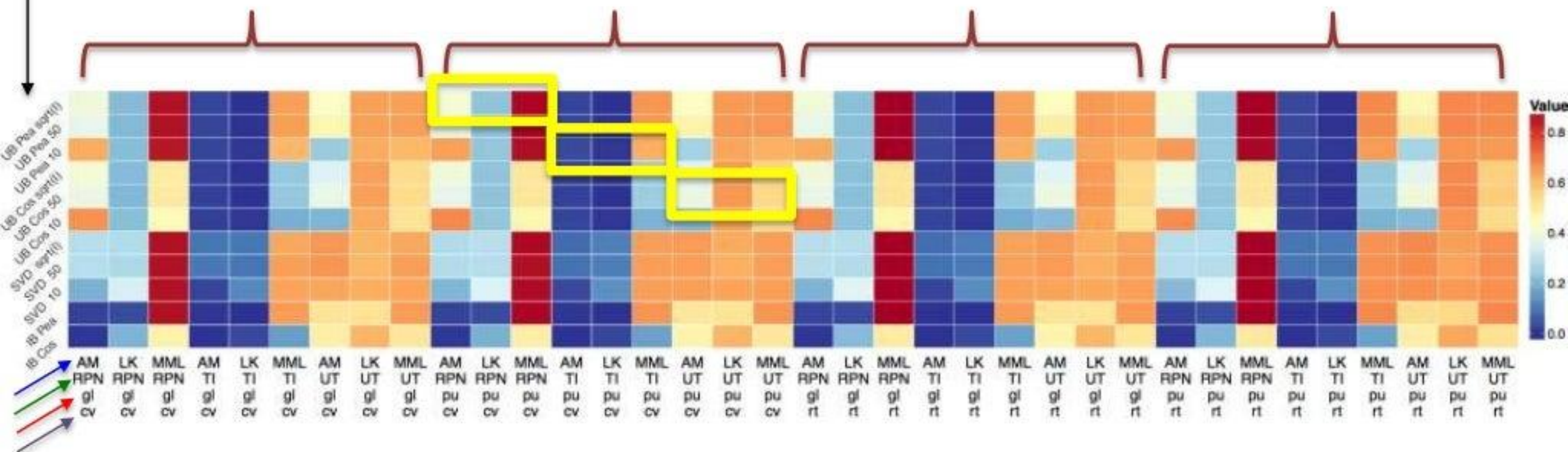
Lenskit vs. Mahout vs. MyMediaLite

Movielens 100k (additional datasets in the paper)

# **AN OBJECTIVE BENCHMARK**

# nDCG@10

Algorithms



## The Frameworks

AM: Apache Mahout  
LK: Lenskit  
MML: MyMediaLite

## The Candidate Items

RPN: Relevant + N [Koren, KDD 2008]  
TI: TrainItems  
UT: UserTest

## Split Point

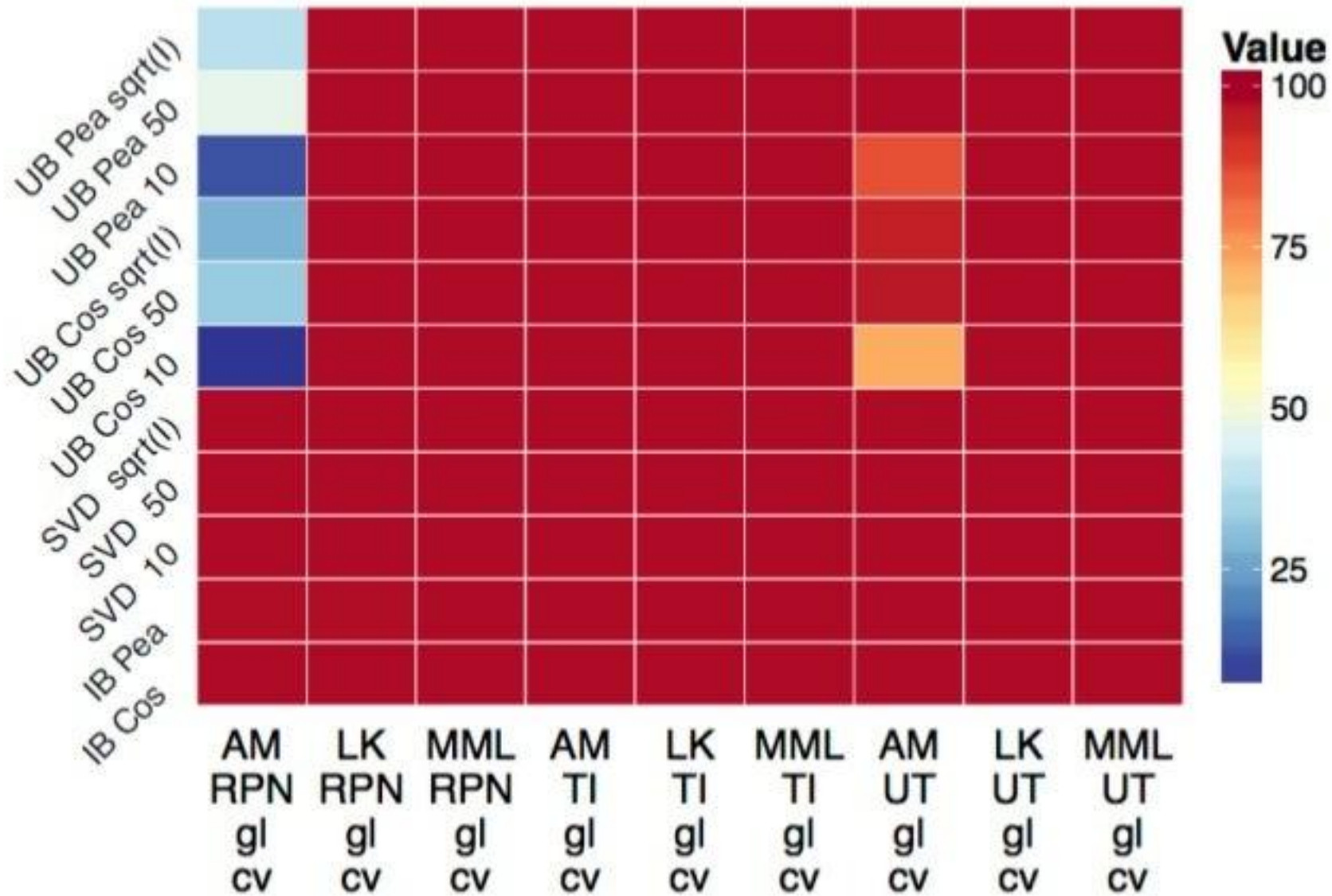
gl: Global  
pu: Per-user

## Split Strategy

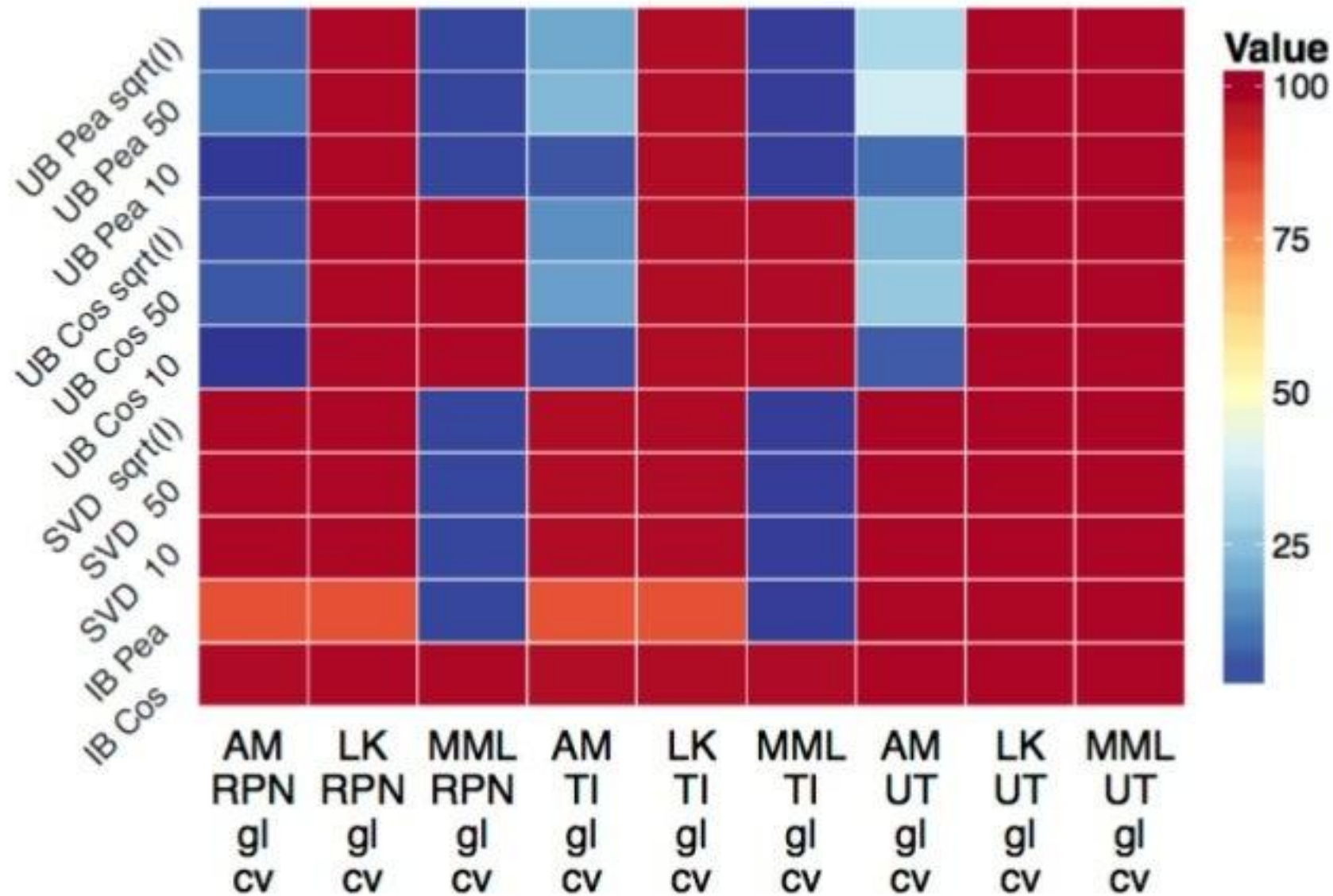
cv: 5-fold cross-validation  
rt: 80-20 random ratio



# User Coverage



# Catalog Coverage





# Time

