

Overview of the 2019 Open-Source IR Replicability Challenge (OSIRRC 2019)

Ryan Clancy, Nicola Ferro, Claudia Hauff,
Jimmy Lin, Tetsuya Sakai, Ze Zhong Wu

Vision



Vision

The ultimate candy store for information retrieval researchers!



Vision

The ultimate candy store for information retrieval researchers!

See a result you like?
Click a button to recreate those results!

Really, any result?
(not quite... let's start with batch ad hoc retrieval experiments on standard test collections)

What is this, really?

Repeatability: you can recreate your own results again

We get this “for free”

Replicability: others can recreate your results (with your code)

Our focus

Reproducibility: others can recreate your results (with code they rewrite)

Stepping stone...

Why is this important?

Good science

Sustained cumulative progress

Armstrong et al. (CIKM 2009): Little empirical
progress made from 1998 to 2009

Why? researchers compare against weak baselines

Yang et al. (SIGIR 2019): Researchers *still*
compare against weak baselines

How do we get there? Open-Source Code!

A good start, but far from enough...

TREC 2015 “Open Runs”
79 submitted runs...



Number of runs successfully replicated

How do we get there?

Open-Source Code!

A good start, but far from enough...

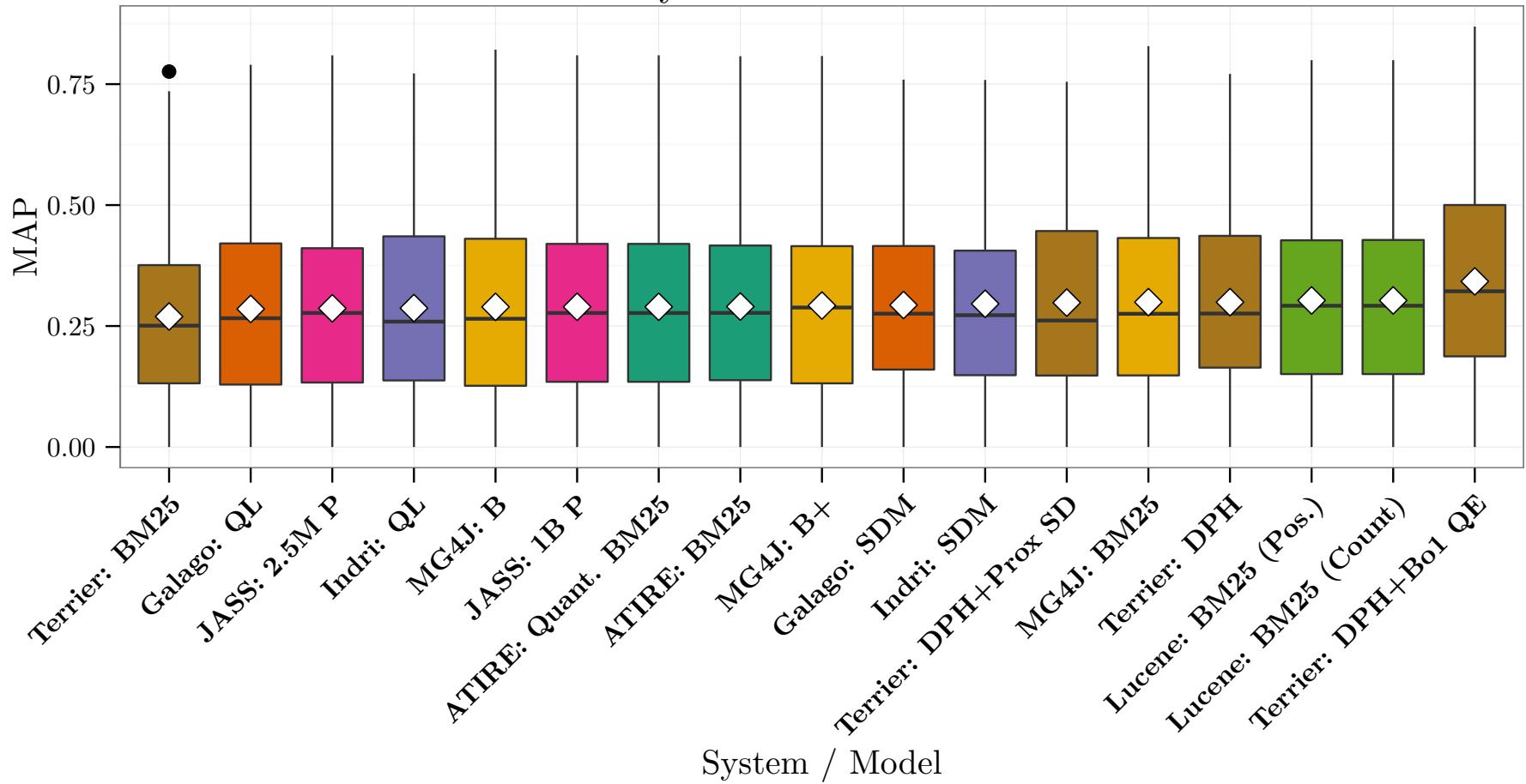
Ask developers to show us how!

Open-Source IR Reproducibility Challenge (OSIRRC),
SIGIR 2015 Workshop on Reproducibility, Inexplicability, and
Generalizability of Results (RIGOR)

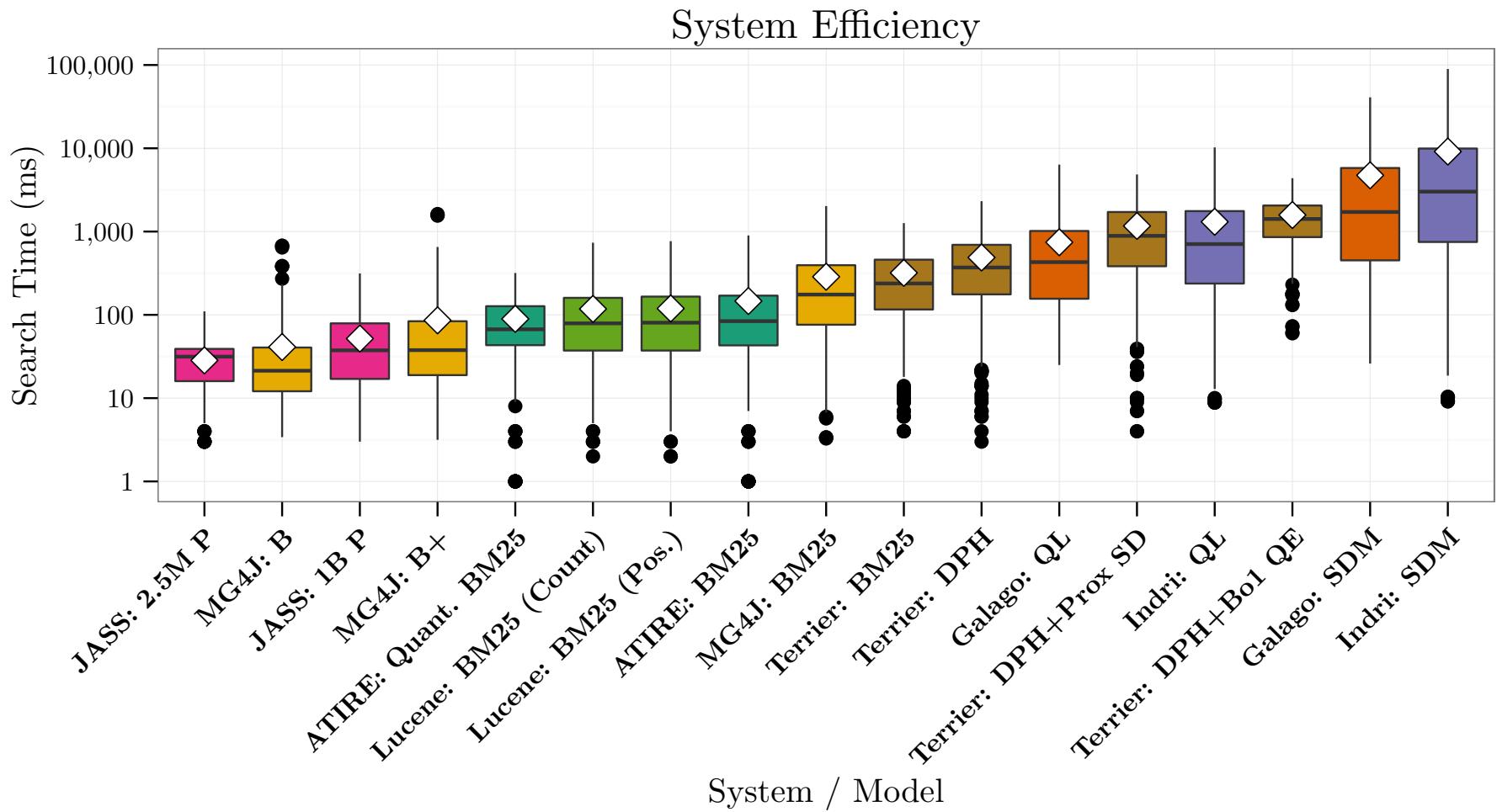
Participants contributed end-to-end scripts for
replicating *ad hoc* retrieval experiments

Lin et al. Toward Reproducible Baselines: The Open-Source
IR Reproducibility Challenge. ECIR 2016.

System Effectiveness

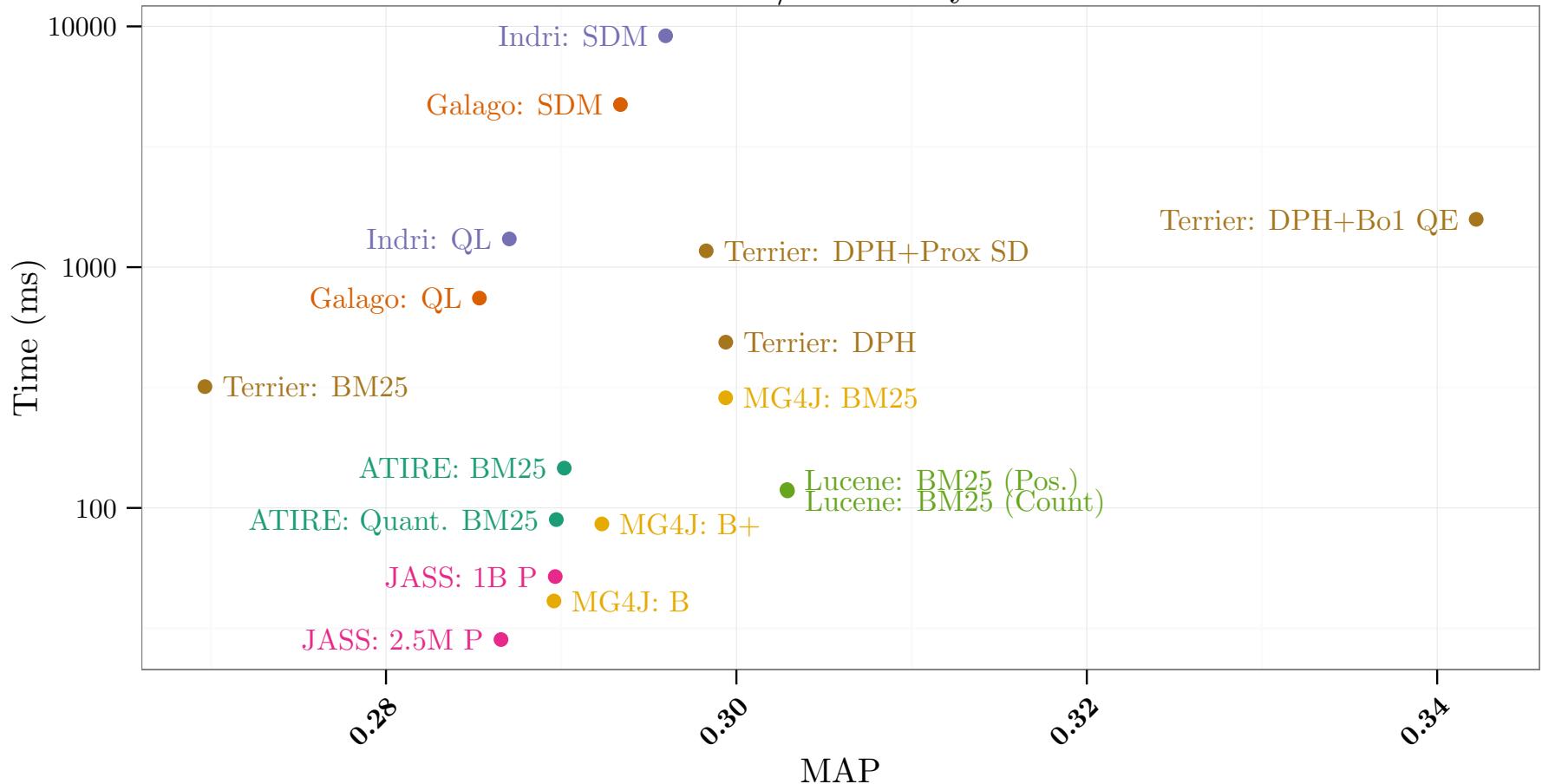


7 participating systems, GOV2 collection



7 participating systems, GOV2 collection

Effectiveness/Efficiency Tradeoff



7 participating systems, GOV2 collection

How do we get there?

Open-Source Code!

A good start, but far from enough...

Ask developers to show us how!

It worked, but...

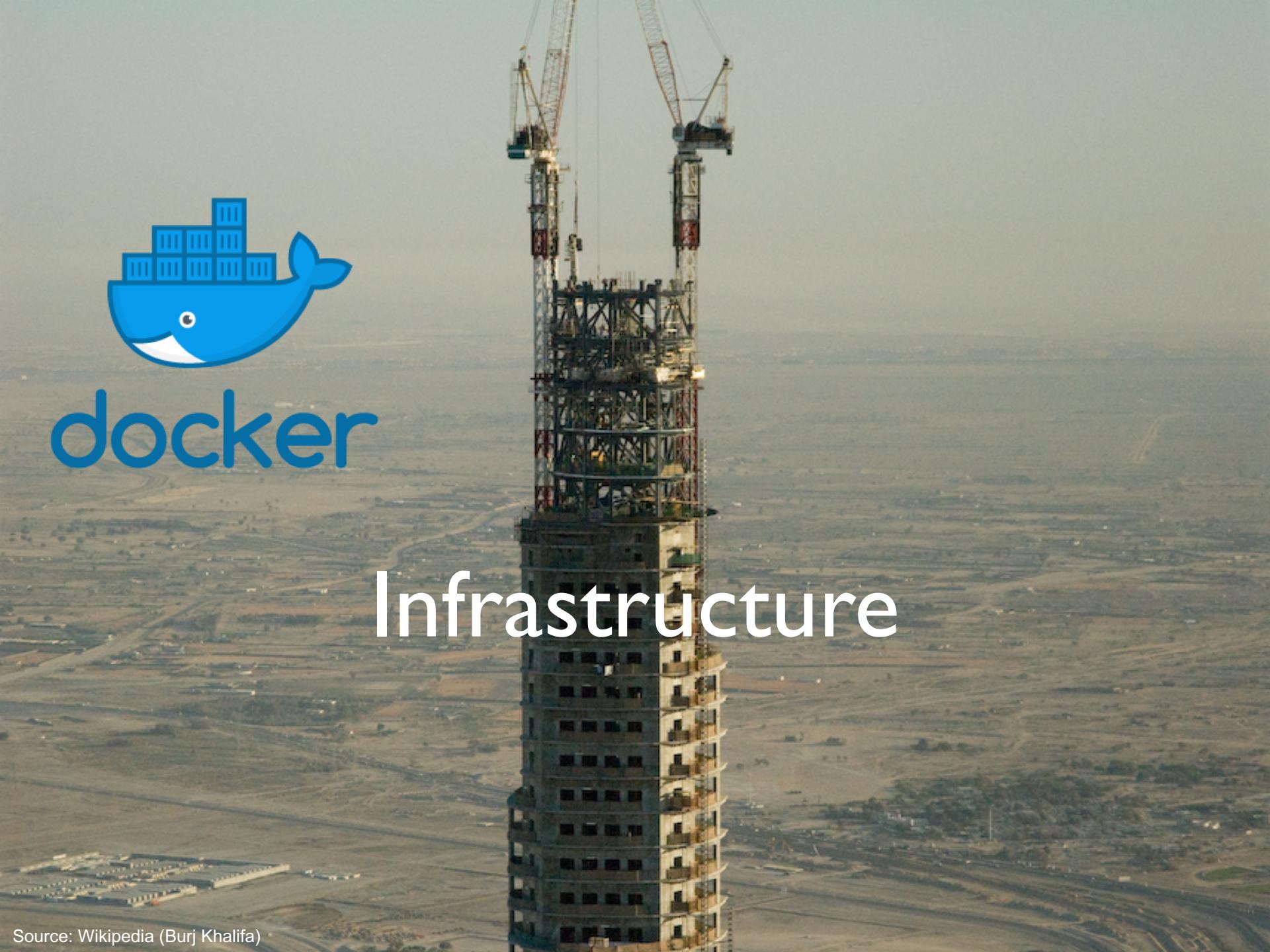
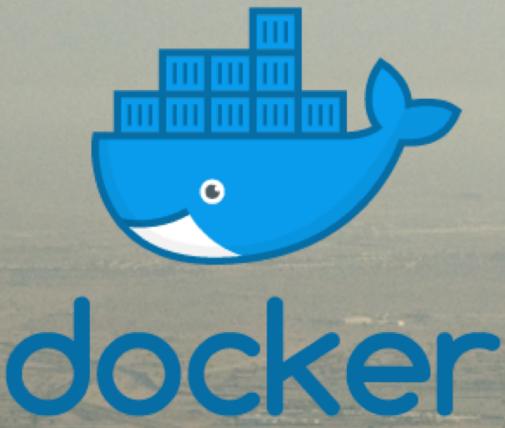
What worked well?

We actually pulled it off!

What didn't work well?

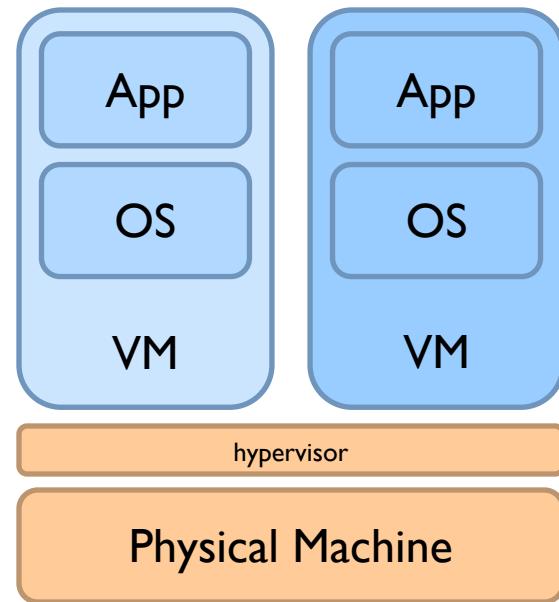
Technical infrastructure was brittle

Replication scripts too under-constrained

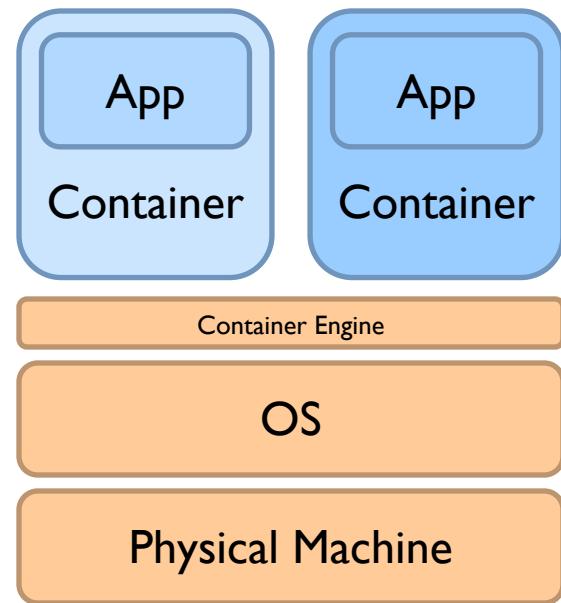
An aerial photograph of the Burj Khalifa during its construction phase. The tower is a massive, multi-layered structure with a complex steel lattice framework. Two large construction cranes are visible at the top of the tower. The surrounding landscape is a vast, flat desert area with some distant structures and roads.

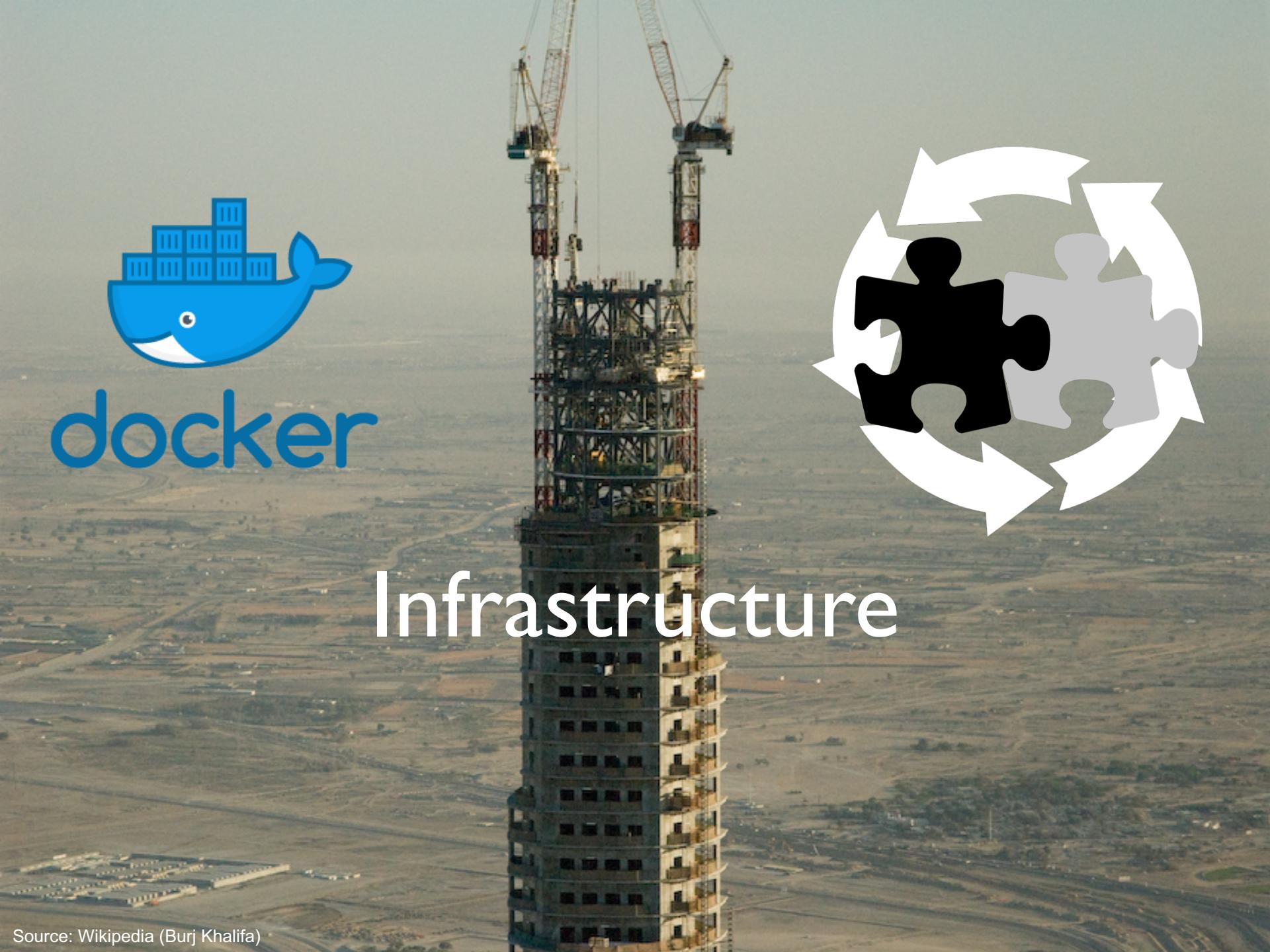
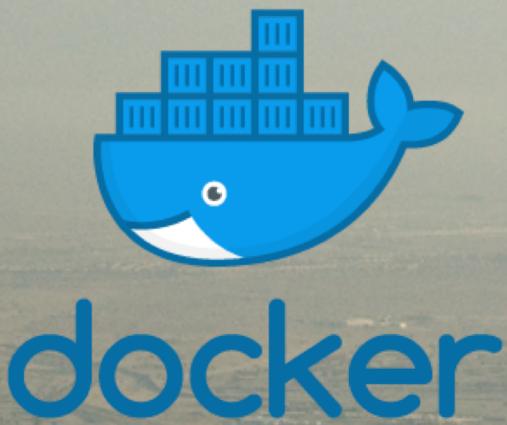
Infrastructure

VMs

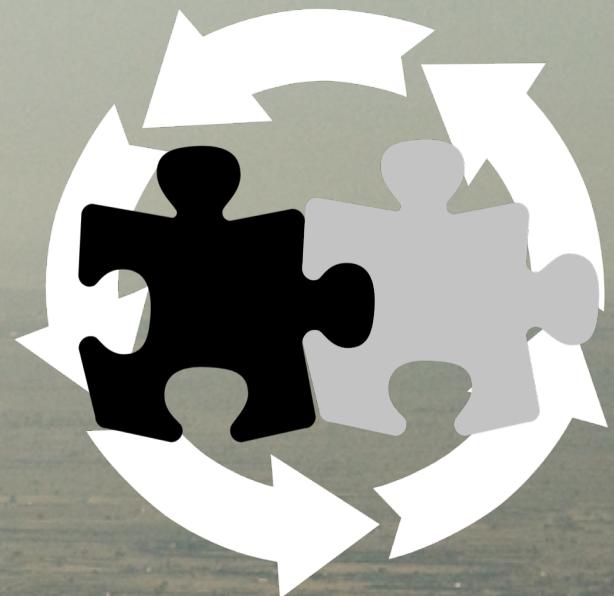


Containers



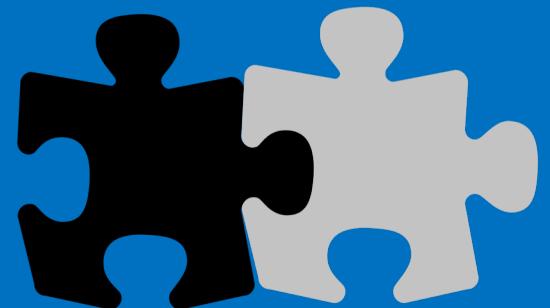
An aerial photograph of the Burj Khalifa during its construction phase. The tower is a tall, slender skyscraper with a complex steel lattice structure. Two large construction cranes are positioned at the top of the tower. The surrounding area is a vast, flat landscape with some smaller buildings and roads.

Infrastructure

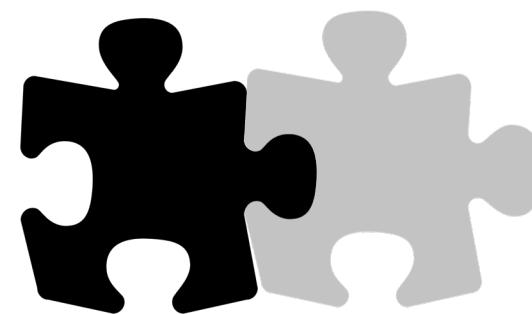
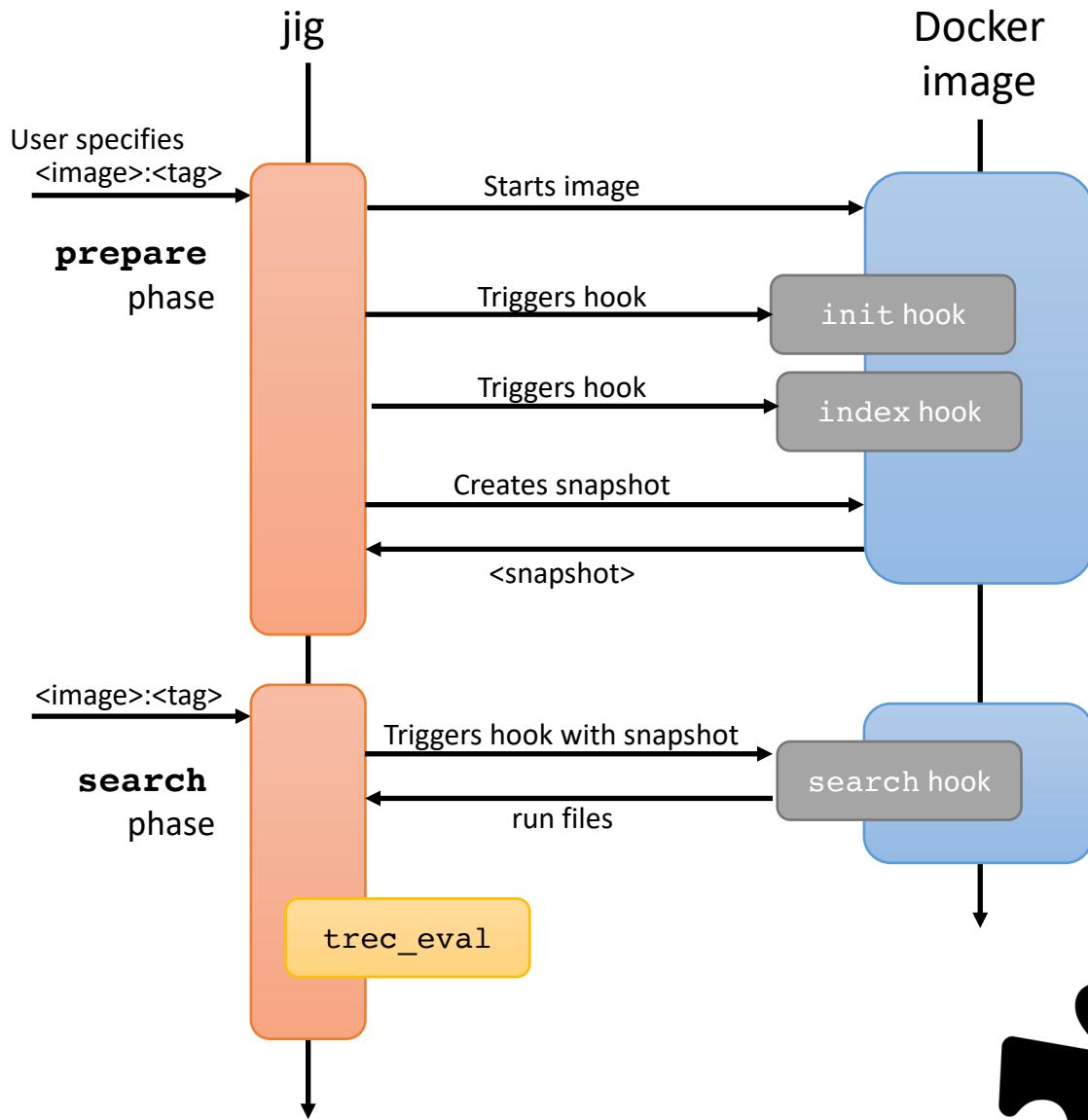


Workshop Goals

1. Develop common Docker specification for capturing *ad hoc* retrieval experiments – the “jig”.
2. Build a library of curated images that work with the jig.



3. Take over the world!
(encourage adoption, broaden to other tasks, etc.)





17 images
13 different teams

Focus on newswire collections: Robust04, Core17, Core18
Official runs on Microsoft Azure

Thanks Microsoft
for free credits!

Anserini (University of Waterloo)
Anserini-bm25prf (Waseda University)
ATIRE (University of Otago)
Birch (University of Waterloo)
Elastirini (University of Waterloo)
EntityRetrieval (Ryerson University)
Galago (University of Massachusetts)
ielab (University of Queensland)
Indri (TU Delft)

IRC-CENTRE2019 (Technische Hochschule Köln)
JASS (University of Otago)
JASSv2 (University of Otago)
NVSM (University of Padua)
OldDog (Radboud University)
PISA (New York University and RMIT University)
Solrini (University of Waterloo)
Terrier (TU Delft and University of Glasgow)

Robust04

49 runs from 13 images

Images captured diverse models:
query expansion and relevance feedback
conjunctive and efficiency-oriented query processing
neural ranking models

Core I7

I2 runs from 6 images

Core I8
I9 runs from 4 images

Robust04

49 runs from 13 images

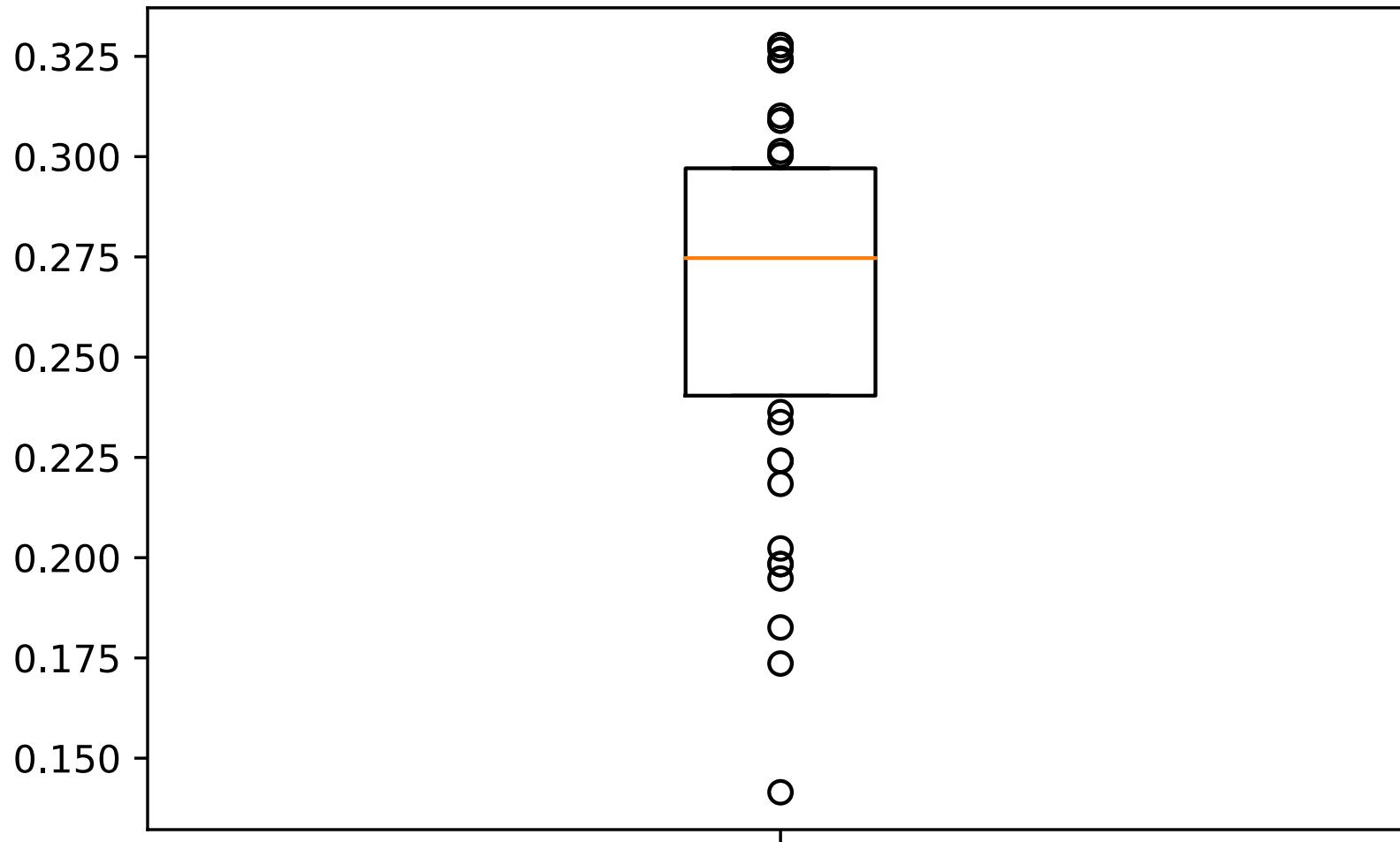


Who won?

A photograph of Senator Bernie Sanders speaking at a podium. He is wearing a dark suit, a light blue shirt, and a dark tie with white polka dots. He has white hair and is wearing glasses. He is gesturing with his right hand while speaking. A microphone is attached to his lapel.

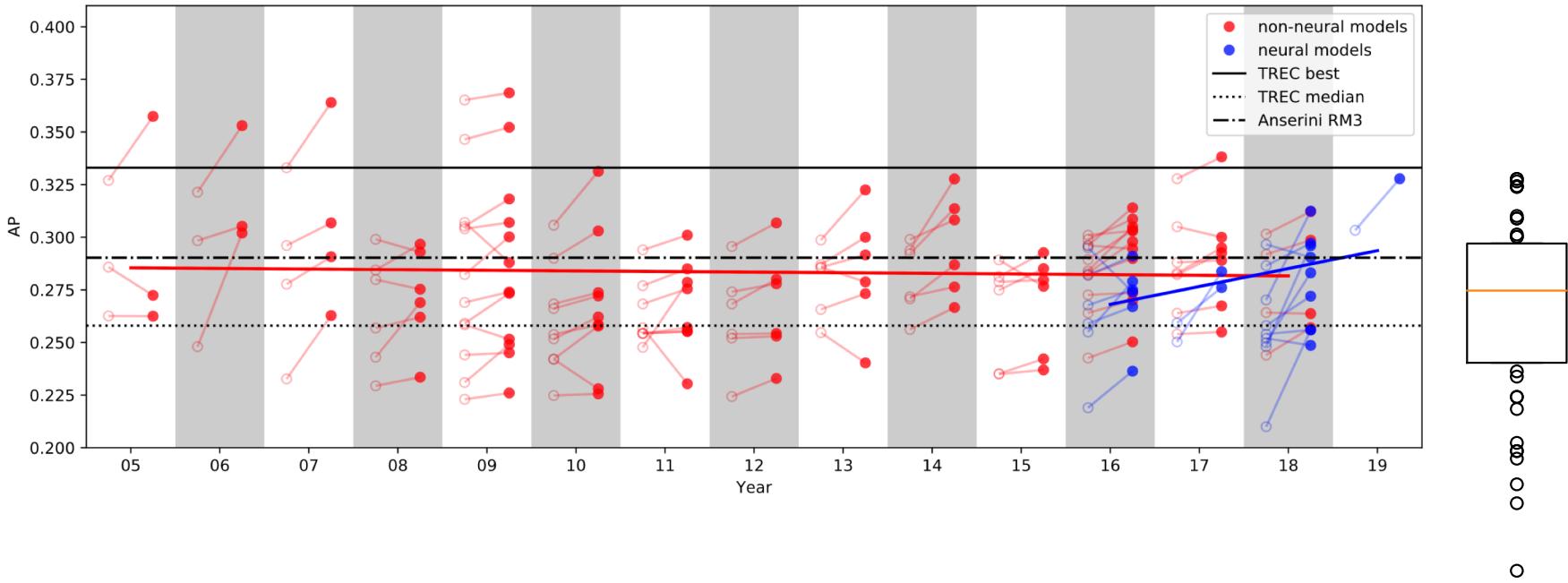
But it's not a
competition!

Robust04 AP Distribution



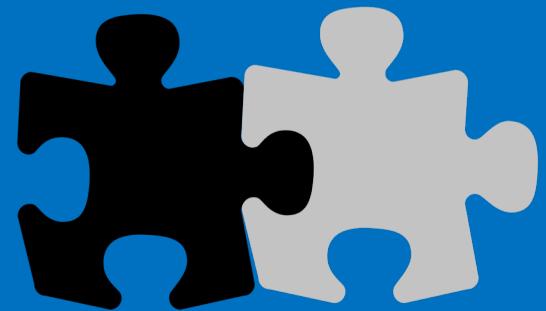
TREC best – 0.333

TREC median (title) – 0.258



Workshop Goals

- ✓ 1. Develop common Docker specification for capturing *ad hoc* retrieval experiments – the “jig”.
- ✓ 2. Build a library of curate images that work with the jig.



- ? 3. Take over the world!
(encourage adoption, broaden to other tasks, etc.)

A photograph of a paved path through tall, golden-brown grass. The path leads towards a range of hills under a cloudy sky.

What's next?