

Big Data

and

Machine Learning

Spring 2023

Hi, I am Panda.

This is a seminar on Big Data and ML

↳ But we are not really going to look at this.

→ Why? Too broad, unclear what I can say that is not better said by an ML class.

Instead going to take a page from last year

Focus on a single area/topic.

Last year: scheduling

This year: Distributed tracing for debugging & profiling.

Three initial questions

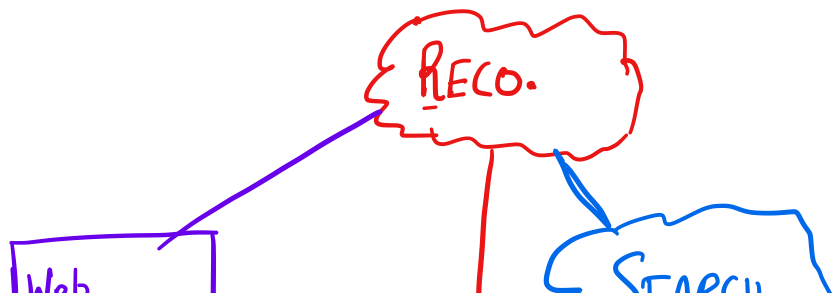
- What?

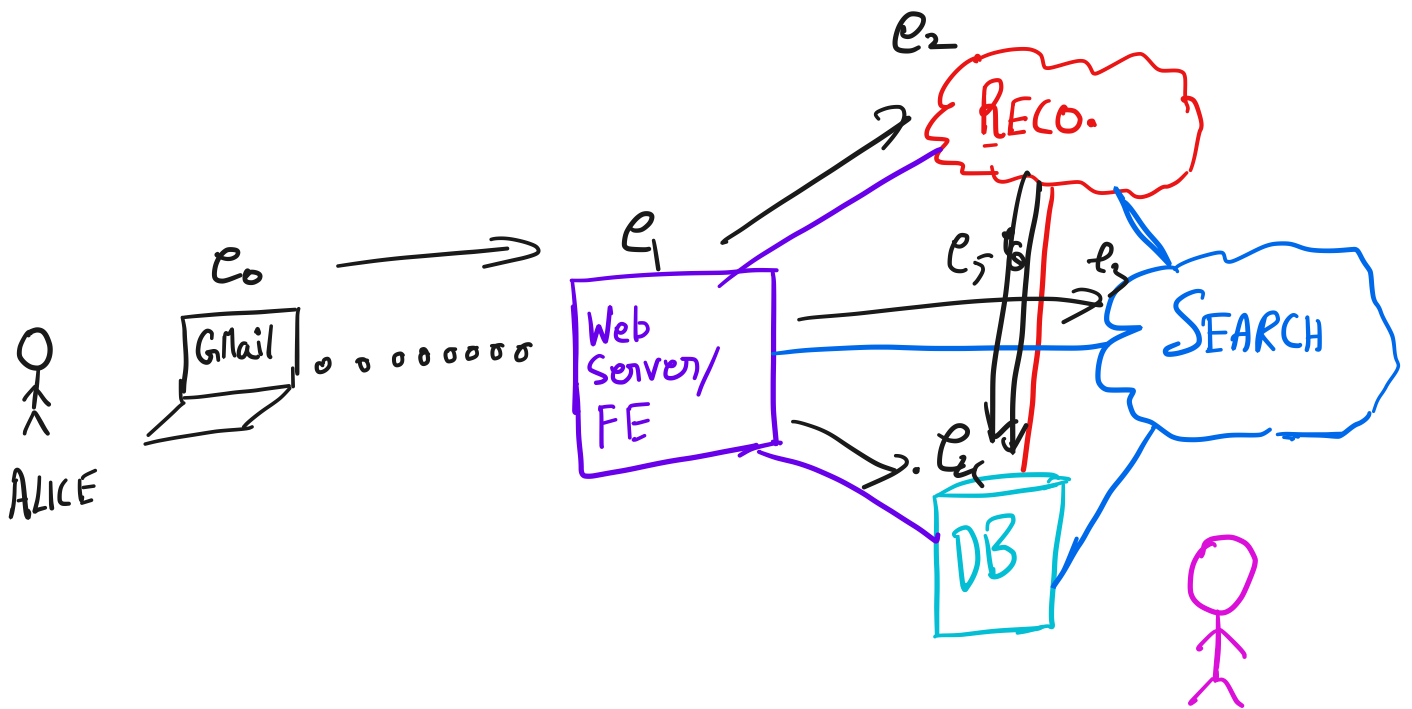
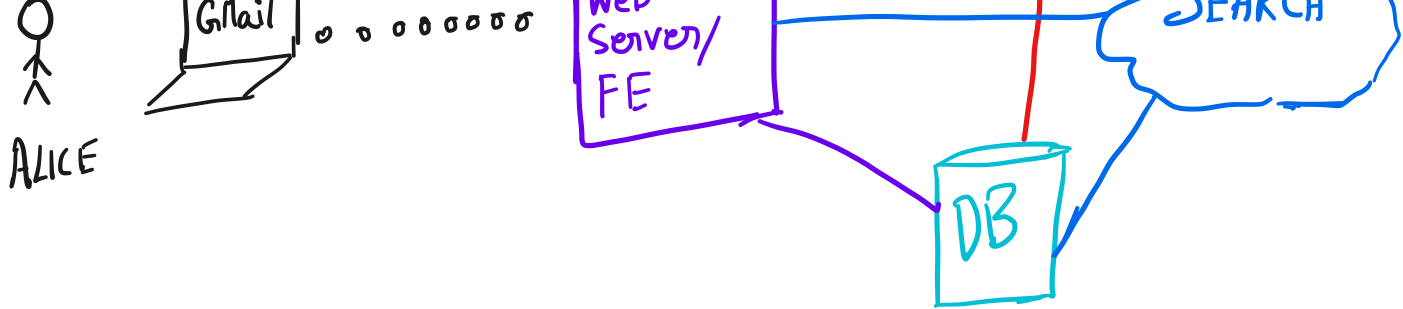
- Why?

- How?

↳ Class mechanics

What?

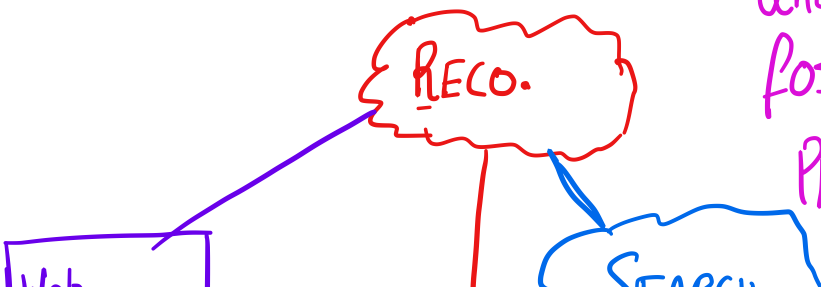


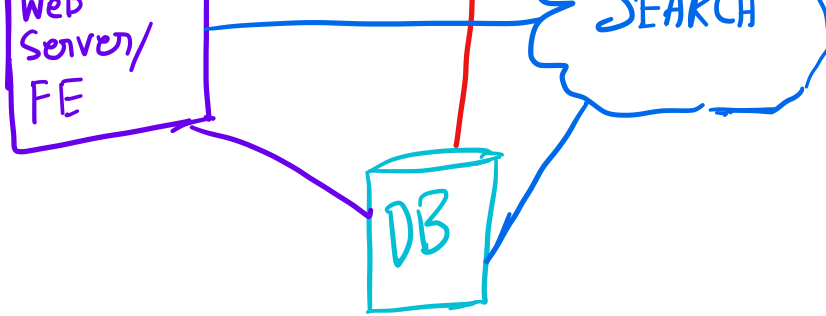


Why is Alice - BETH SRF
unable to search
for e-mails about
PIZZA?

• Why is Alice
unable to search
for e-mails about
PIZZA?

BETH SRF





- Why is search using so many resources?
- Why are recommendations delayed?

- OFTEN THESE PROBLEMS ARE DUE TO INTERACTION B/W SERVICES

FOR EXAMPLE

- RECOMMENDATION ENGINE ASSUMES LOW DB LATENCY
→ MAKES MANY QUERIES
 - DB LATENCY DEPENDS ON # OF CONCURRENT QUERIES
→ LATENCY INCREASES AS # OF CONCURRENT QUERIES INCREASE
 - A RECENT SEARCH FEATURE INCREASES # OF DB & RECOMMENDATION QUERIES
- ⇒ - SEARCH REQUIRES MORE RESOURCES
- RECOMMENDATIONS ARE SLOW

NEED TOOLS AND METHODS TO REASON ACROSS

ALL OF THESE SERVICES

- EACH OF WHICH MIGHT BE COMPRISED OF SEVERAL PROCESSES / SERVICES

DISTRIBUTED TRACING IS ABOUT COLLATING AND CORRELATING INFORMATION ACROSS PROCESSES

REQUIREMENTS?

- o Record time - response time
Order
- o Low overhead
- o Always up / fault tolerant
- o Log levels

Why?

- USEFUL ACROSS DOMAINS
 - ↳ UNDERSTANDING BOTTLENECKS / PROBLEMS IN ML
 - " " " " IN DATA PROCESSING
 - . . .
- WE CAN ACTUALLY GENERATE TRACES, USE & ANALYZE THEM
 - ↳ TASK FOR NEXT CLASS.
 - DOESN'T REQUIRE ACCESS TO FANCY H/W, ETC.
 - WE HAVE PRETTY REASONABLE VISIBILITY INTO WHAT PEOPLE DO IN PRACTICE.
- THERE CONTINUE TO BE MANY OPEN QUESTIONS & PROBLEMS THAT NEED TO BE SOLVED.

How?

- o READ & LEARN FROM
 - ↳ A RECENT BOOK FROM PEOPLE AT MICROSOFT, TWITTER, ETC. ON THE TOPIC
 - ↳ ACCESS ONLINE, FOR FREE, FROM NYU LIBRARY

→ PAPERS

→ DOCUMENTATION & CODE

○ GET YOUR HANDS DIRTY WITH COLLECTING & ANALYZING TRACES

→ NO STENCIL CODE -
REALLY WRITE SOMETHING THAT
EXHIBITS PHENOMENON YOU THINK ARE
INTERESTING.

→ THIS WEEK ○ GET STARTED WITH GENERATING TRACES
FROM A SINGLE PROCES

↳ EVERY FEW WEEKS SNAPSHOT EXERCISES & SUBMIT
AS HW (3 × 10%)

- GENERATING & COLLECTING TRACES

- VISUALIZING & ANALYZING TRACES

- USING TRACES TO DEBUG PROBLEMS

○ FINAL PROJECT (25%)

- INTERESTING WAYS TO COLLECT OR USE TRACES

TRACE FRAMEWORKS

- EXTENSIONS to EXISTING TRACE

- ...

o EXAMS (10+20)

↳ UNDERSTANDING OF CLASS MATERIAL

o PARTICIPATION

↳ IN CLASS (5%)

→ HELPING PEOPLE ONLINE, POSTING TUTORIALS OR INSTRUCTIONS, PRESENTING TOOLS/TECHNIQUES IN

CLASS (10%)

WHAT GOES INTO A TRACE

REQUIREMENTS

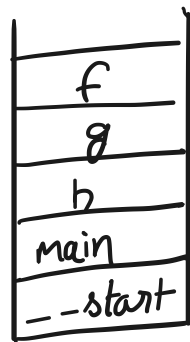
o TRACK CAUSALITY

WHAT?

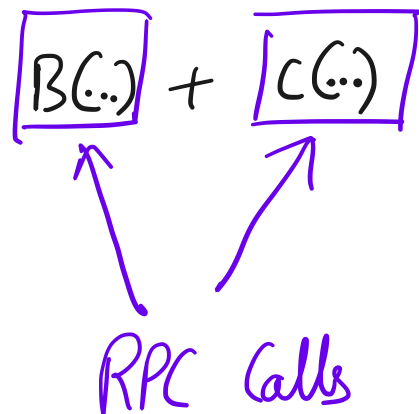
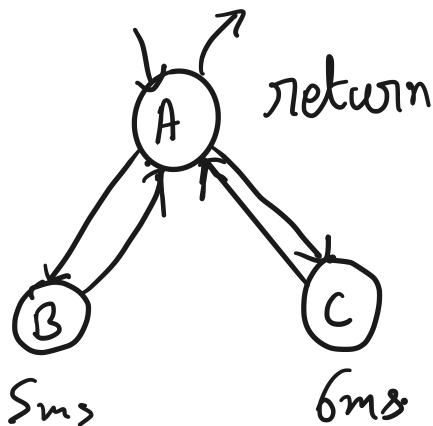
Why?

o SUPPORT CONCURRENT CALLS

STACK TRACES FOR PROGRAMS



'bt' in gdb



EXECUTION TIME

EXECUTION TIME

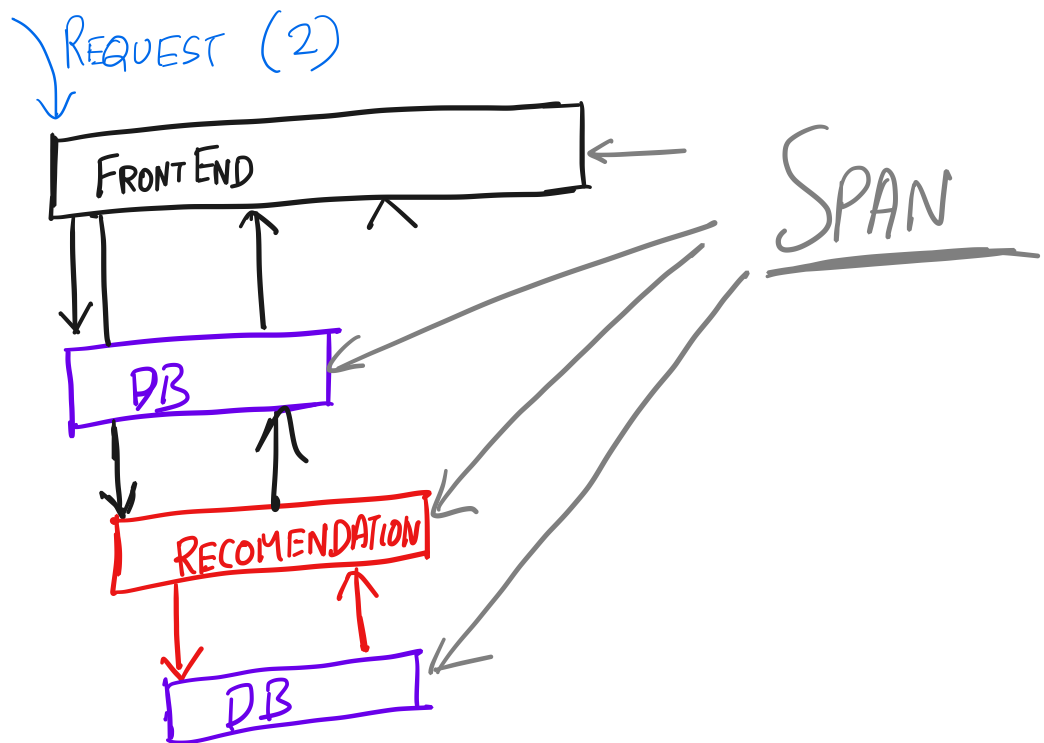
- o TRACK TIME (& OTHER ATTRIBUTES)

THE STRUCTURE OF TRACES

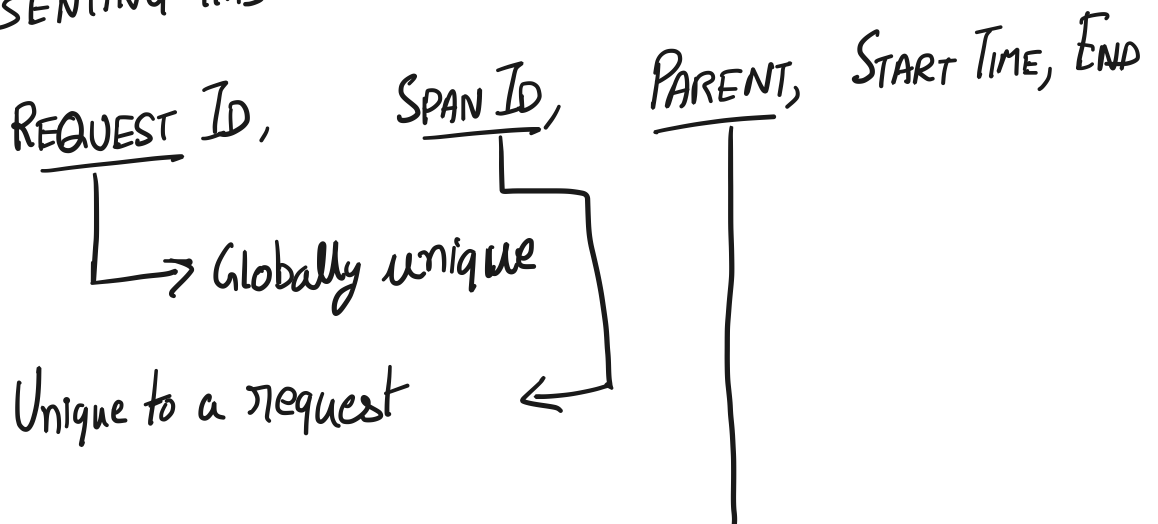
- o USE IDs (Request ID or Trace ID) to DISTINGUISH BETWEEN REQUESTS.


WHAT IS A REQUEST?

◦ WITHIN A REQUEST ◦ MODEL CAUSALITY AS A TREE.



◦ REPRESENTING THIS



Span ID on \perp 

o WHERE ARE THESE STORED?

o HOW COLLECTED?

o OVERHEADS