

PORTFOLIO ANALYTICS IS A BIG DATA PROBLEM

IS IT?

ALL AND ONLY MY PERSONAL VIEWS HERE

2/5/2019

Luca S.V. Spampinato – Bloomberg L.P.

lspampinato1@bloomberg.net

ALWAYS START WITH A QUOTATION

Yngve Slyngstad, since 2008 CEO of the world's biggest sovereign wealth fund, Norges Bank Investment Management (> \$1Trillion AUM)

Bloomberg News Interview 02/02/2019

BM: What do you think is different today in terms of managing money vs. how it was when you started off at the fund in 1998 or even earlier at Storebrand?

YS: It is still about information processing, but the amount of information that is available is of course increasing every year, and the frequency of that information is just getting faster and faster. You have to cut through that and find what is essential. With this kind of a skill set it is very difficult to see who has got it and who hasn't got that ability, but I think it's one way of distinguishing.

BM: How is the advance of technology and AI changing investing?

YS: My own guess is that it's going to dramatically change quantitative investing and particularly risk-factor investing, probably more so than traditional active management.

AI IN PORTFOLIO ANALYTICS IS RUNNING ALONG MAINSTREAM

Optimization in various flavours is a best hit

But there are actual, production ready, “side” results:

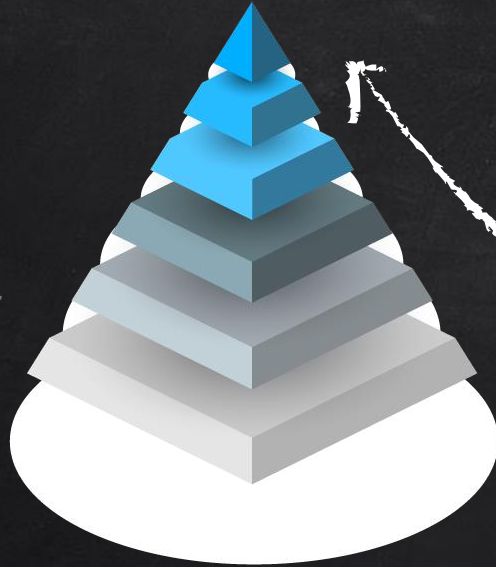
- x Liquidity analysis
- x Forecast of cash flows
(subscriptions, and redemptions in Mutual Funds)

One of the pros* of wide spreading AI in finance is that it brings focus on Data Processing

*(the only?)

THE ANALYTICS PYRAMID

AI
ML
Data Science
Analytics
Data Engineer
Big Data



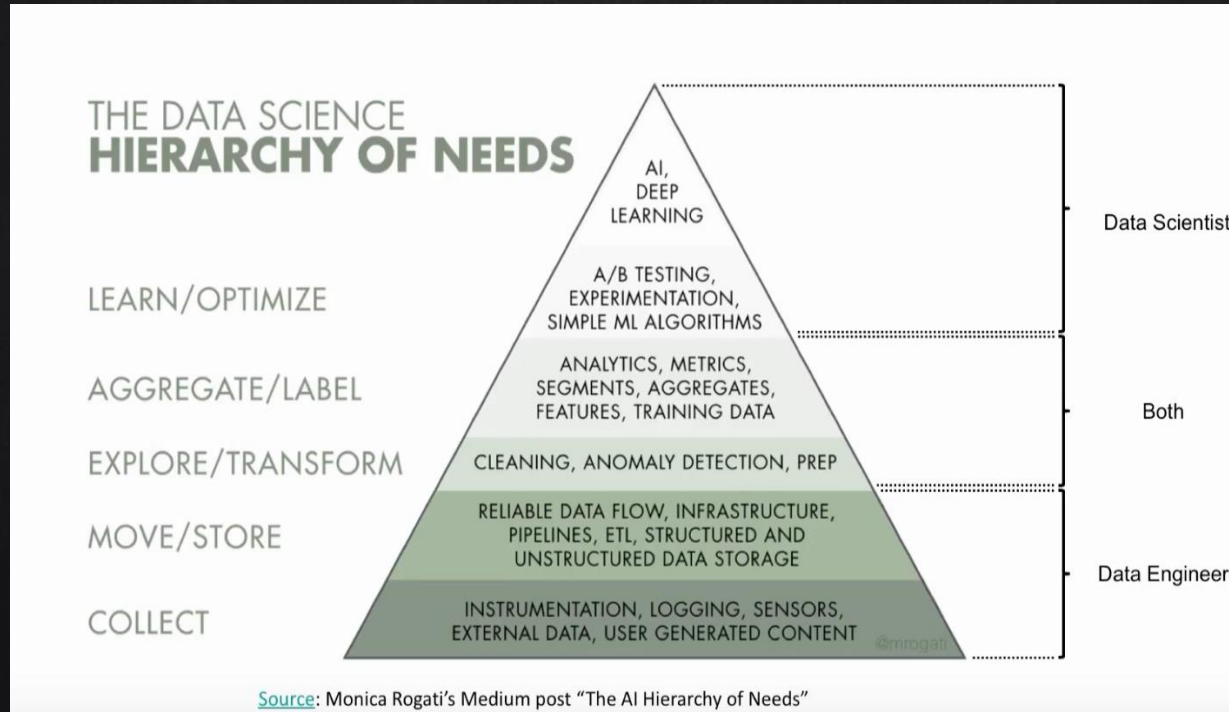
Portfolio Analytics is a Data Management problem anyway:

- X Heterogeneous Data Sources
- X Size of Portfolios
- X Analytics demand

...AND CLIMBING UP THE PYRAMID CANNOT BE AVOIDED

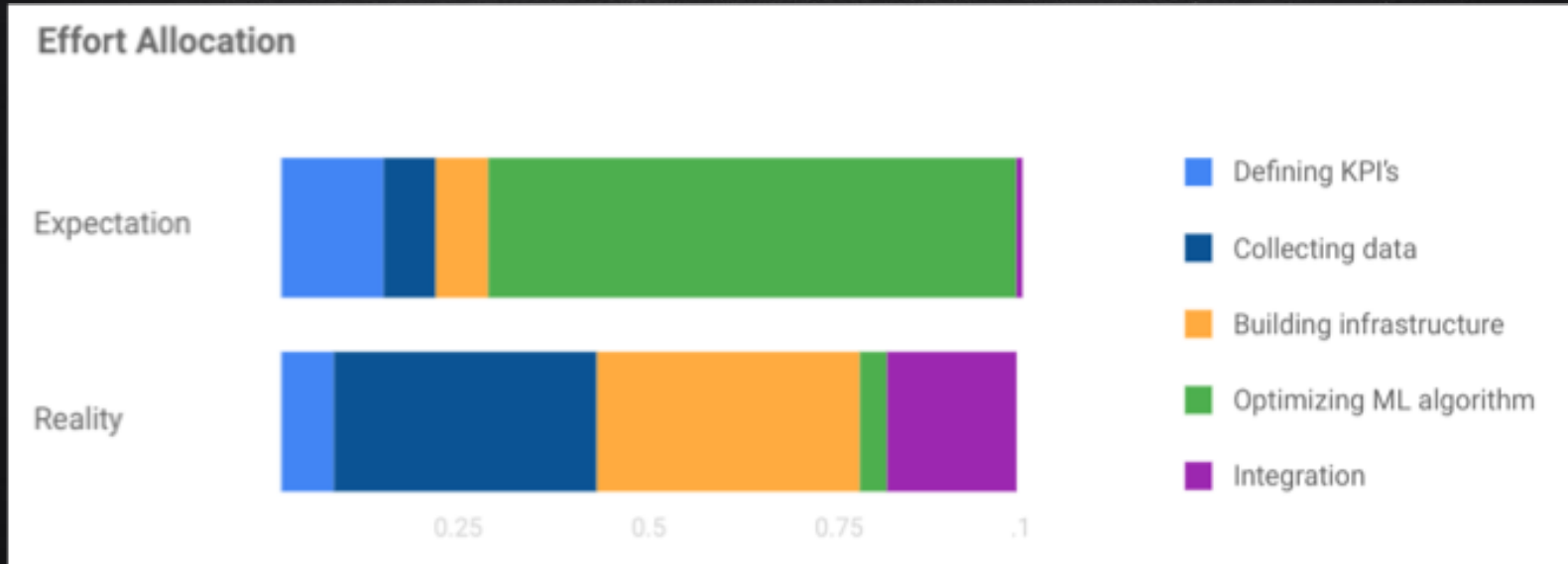
NEW ROLES

and, by the way, Data Scientists cannot be left alone



THE SURPRISE

ML projects' pain, in retrospective



... does this sound familiar?

EXPANDING DATA UNIVERSE 1 – ANALYTICS DEMAND

- X Asset owners sophistication, at any level (e.g. Risk Monitoring is the word in wealth management)
- X Regulation

- X Scenario Analysis (hundreds of thousands of scenarios, full repricing, for VaR, Asset Liability Modeling, etc.)
- X Factor based analysis (thousands of factors)
- X Multi-factor, multi-dimension Performance&Risk Attribution

EXPANDING UNIVERSE 2 - NUMBER, SIZE AND COMPLEXITY OF PORTFOLIOS

NUMBER

- x Wealth management
- x Proliferation of funds

SIZE

- x Asset owners aggregation
- x Custom, structured Benchmarks

COMPLEXITY

- x Pervasive multi-asset
- x Custom classifications vs. Strategies

EXPANDING UNIVERSE 3 - DATA QUANTITY AND COMPLEXITY

HOLDINGS

- X Decades
- X daily

PRICES:

- X Many sources and formats
- X different frequencies
- X "round-the-clock" alignment

REFERENCE:

- X Terms & Conditions
- X All asset classes & types
- X UDIs

AUXILIARY:

- X Rates/Curves
- X Classifications
- X Indexes and benchmarks

TRANSFORM:

- X Currency rates
- X C.A.s
- X Distributions & Accruals

BENCH.

- X Decades
- X Daily
- X Multi-family
- X Prices
- X 100,00s

RISK FACTORS FUEL:

- X Fundamentals & Estimates (Fiscal alignment)
- X Ratings
- X Credit Data

BEYOND:

- X News & Events
- X Supply Chain
- X Sector Trends

NUMBERS (NOWADAYS IN THE PORT ECOSYSTEM)

>1MILLION SECURITIES
Not including UDIs

>12,000 FIELDS
Not including Custom data

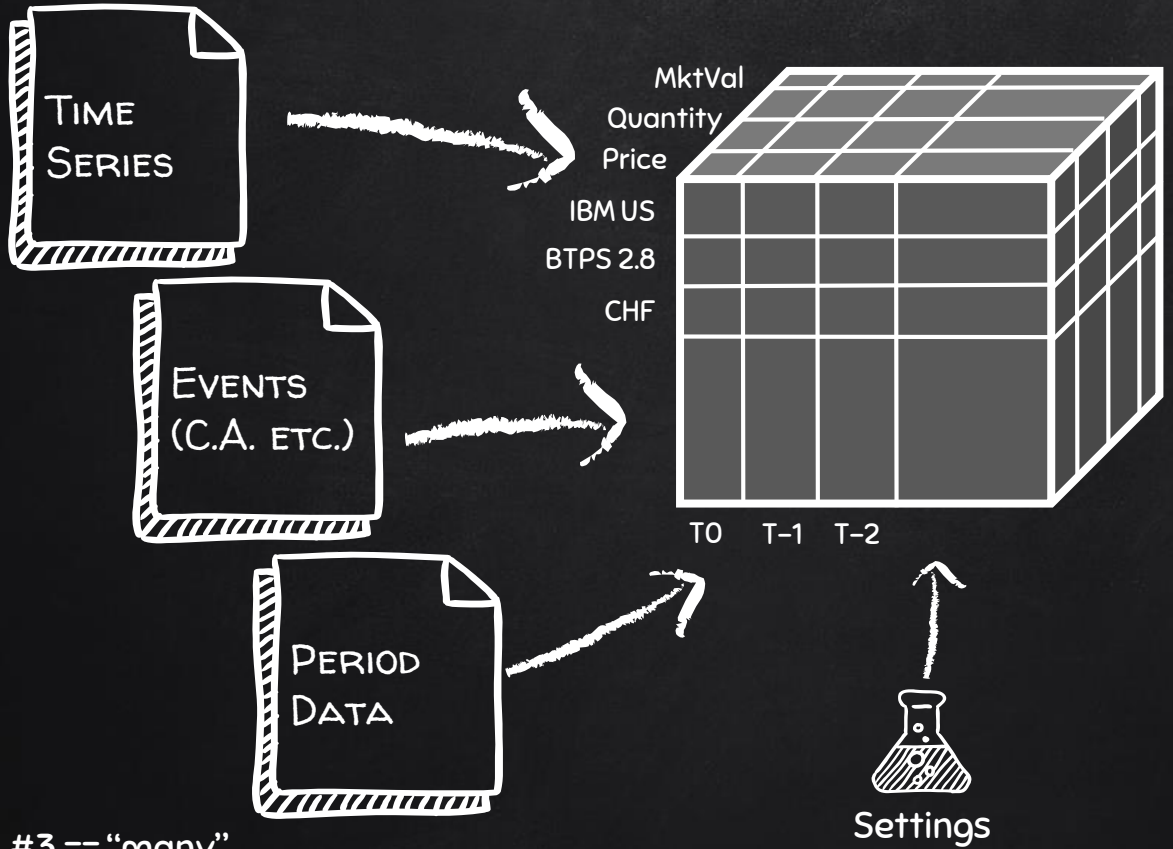
~800 TB DAILY DATA FLOW
i.e. ~24 PB/month

16 UPDATES/DAY

~80,000 AVG. QUERY PER HOUR
160,000 Peak

PORTFOLIO ANALYTICS DATA REQUIREMENTS

BOIL DOWN TO A (BIG) CUBE



Key factors are:

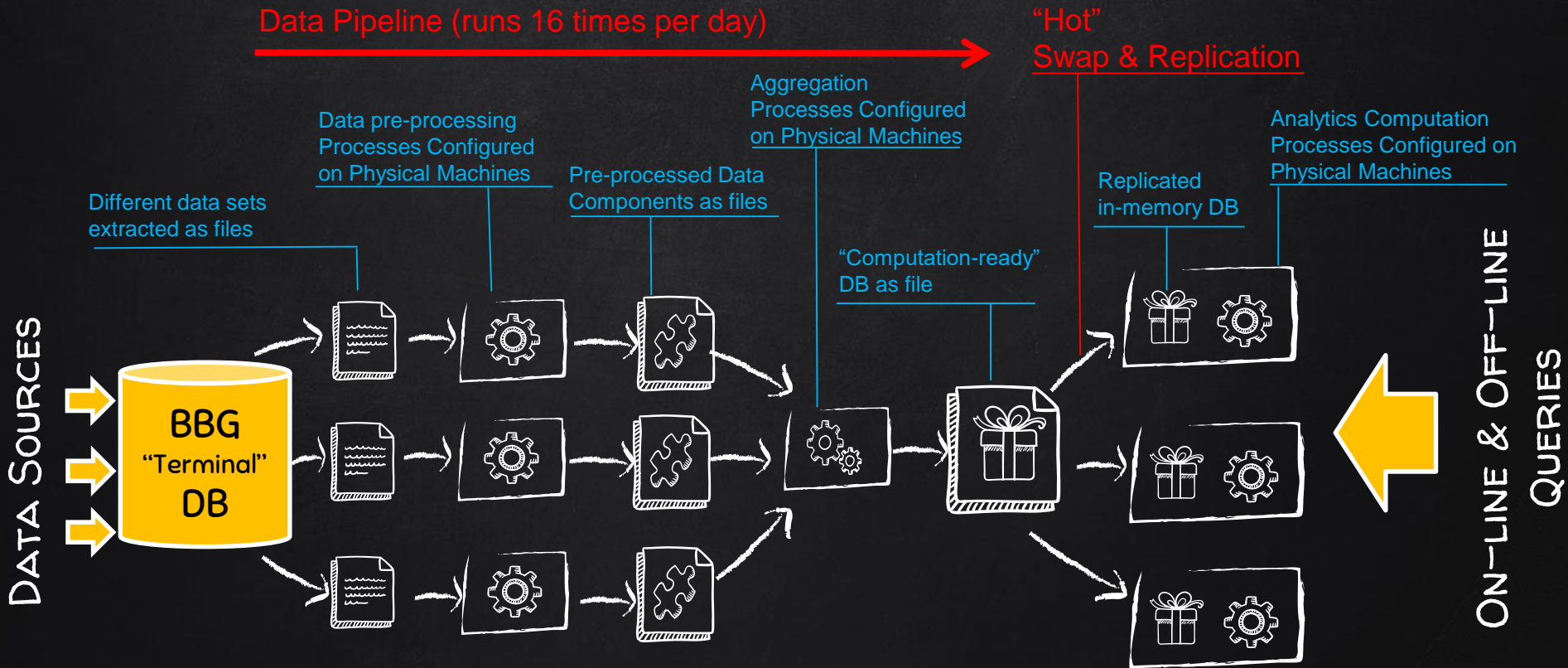
- x Bring the data in the cube
- x Compute in the cube

(easy to say, isn't it?)

#3 == "many"

WHAT WE DO NOW ('90s ON STEROIDS)

A Traditional large scale architecture: scale-up means re-configuration



#3 == "many"

"UNO VALE UNO" I.E. "UNO ALLA VOLTA, PER CARITA"

350+100 GB RAM

| Portfolio | T0 | T-1 | T-2 |
|-----------|----|-----|-----|
| IBM US | | | |
| BTPS 2.8 | | | |
| CHF | | | |



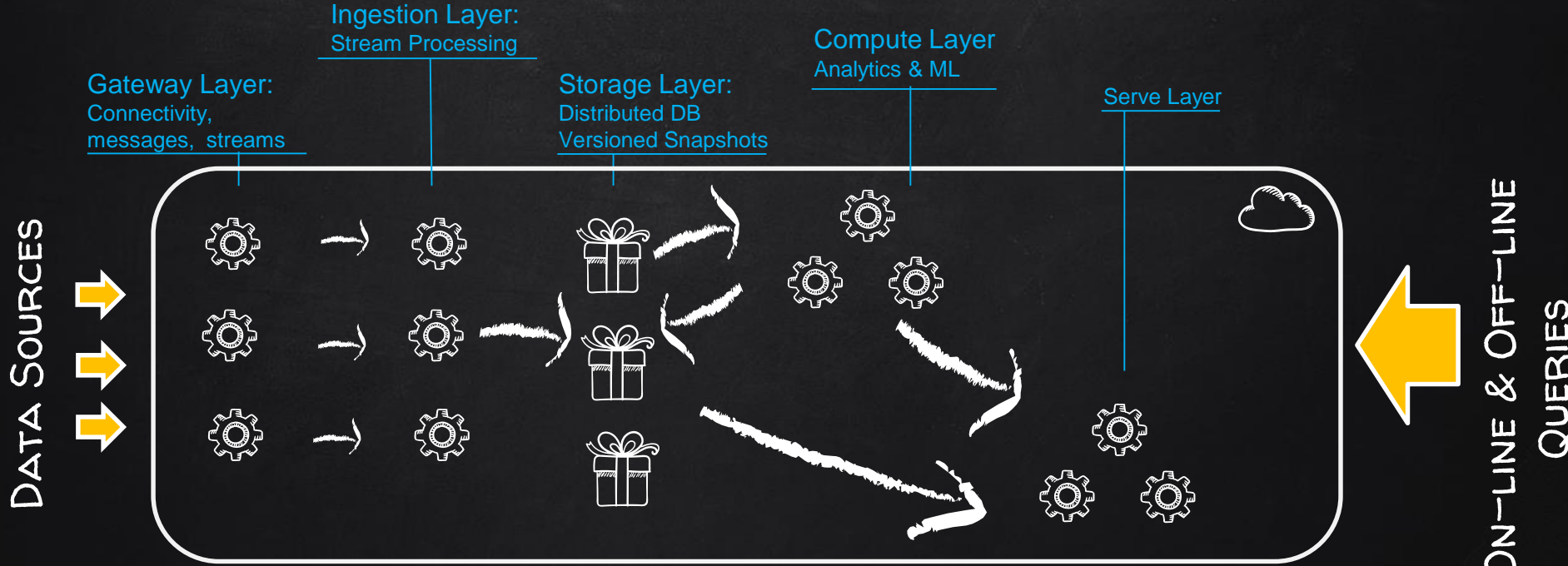
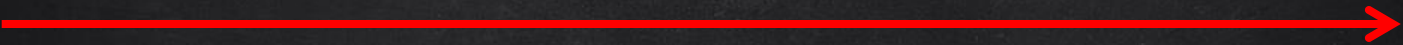
ON-LINE & OFF-LINE
QUERIES

#3 == "many"

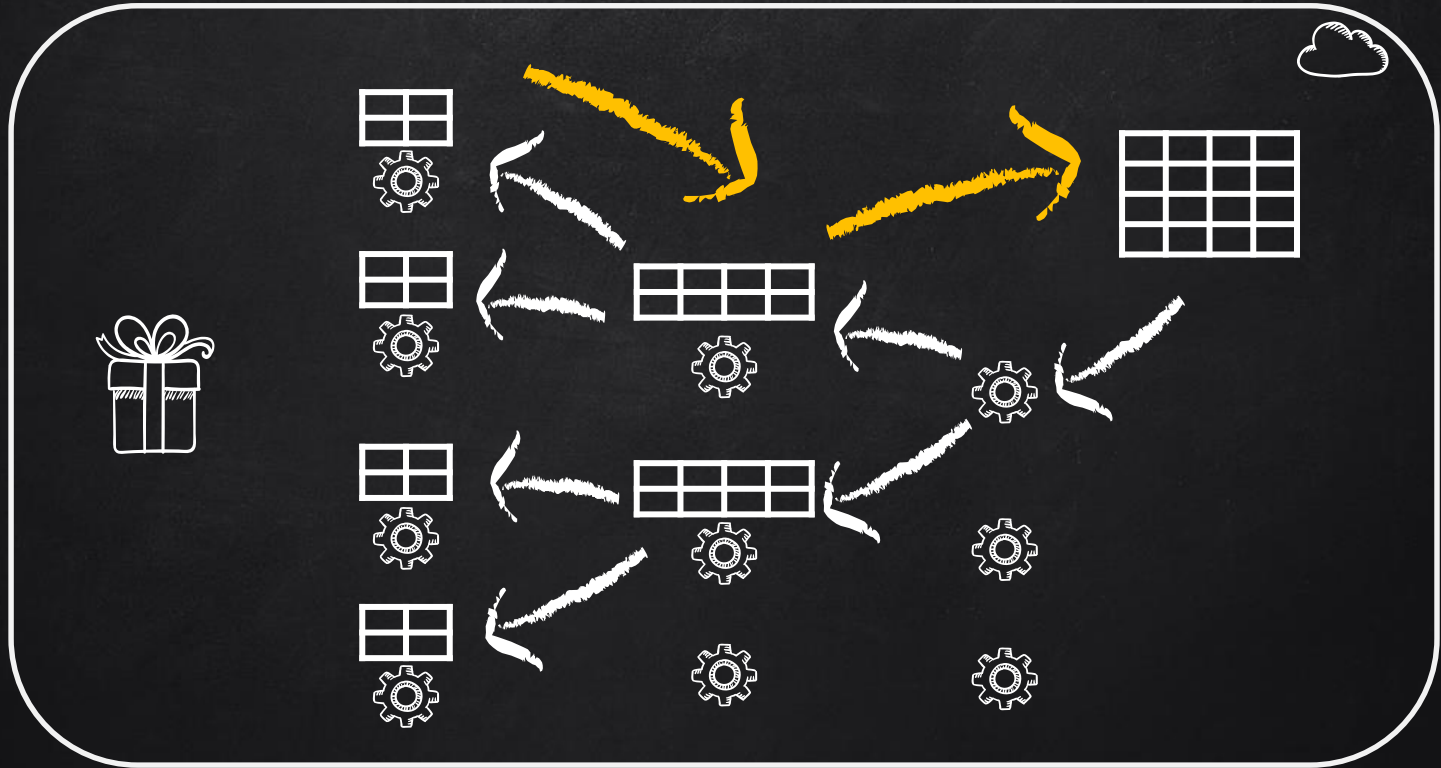
THE TARGET

Virtualized Applications is the key: scale-up is automatic and self-adapting

Continuous run



DISTRIBUTED COMPUTING I: PARADIGMATIC MAP-REDUCE



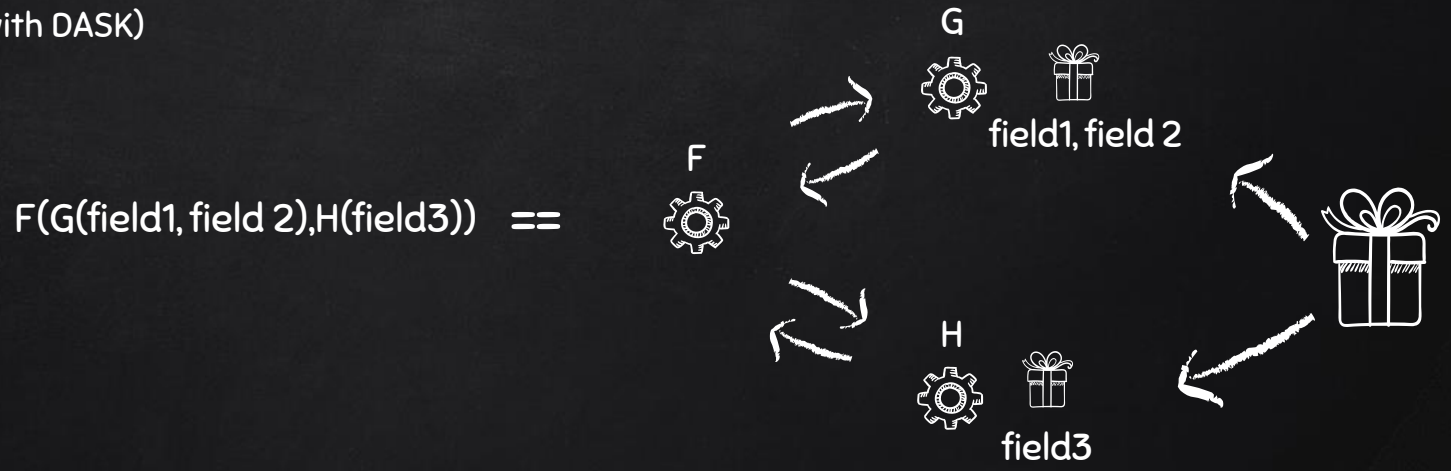
ON-LINE & OFF-LINE
QUERIES

DISTRIBUTED COMPUTING II: TRANSPARENT DISTRIBUTION

Computation Tree is meta-described:

- X Distribution is transparent (independent branches are executed in parallel)
- X Lazy evaluation brings the right data to the right computation nodes

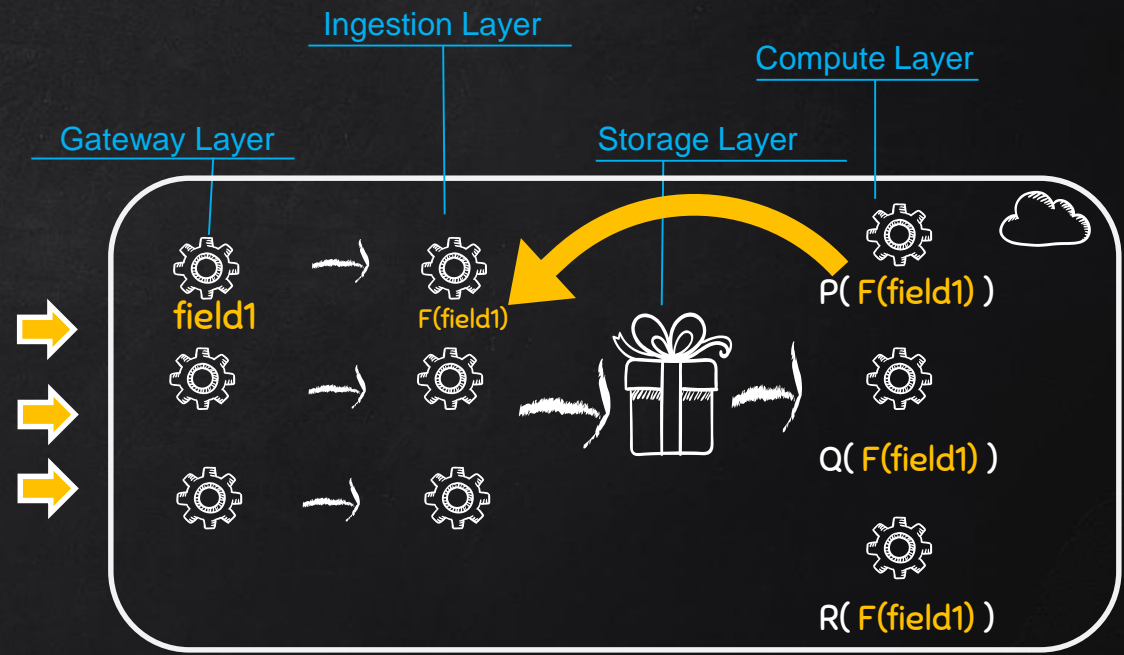
(e.g. Python code with DASK)



DISTRIBUTED COMPUTING III: GREEDY PIPELINE

Data-bound computation modules can be easily* pushed backward in the Data Pipeline

* well, yes, they have to be properly packaged



THANKS & QUESTIONS

lspampinato1@bloomberg.net