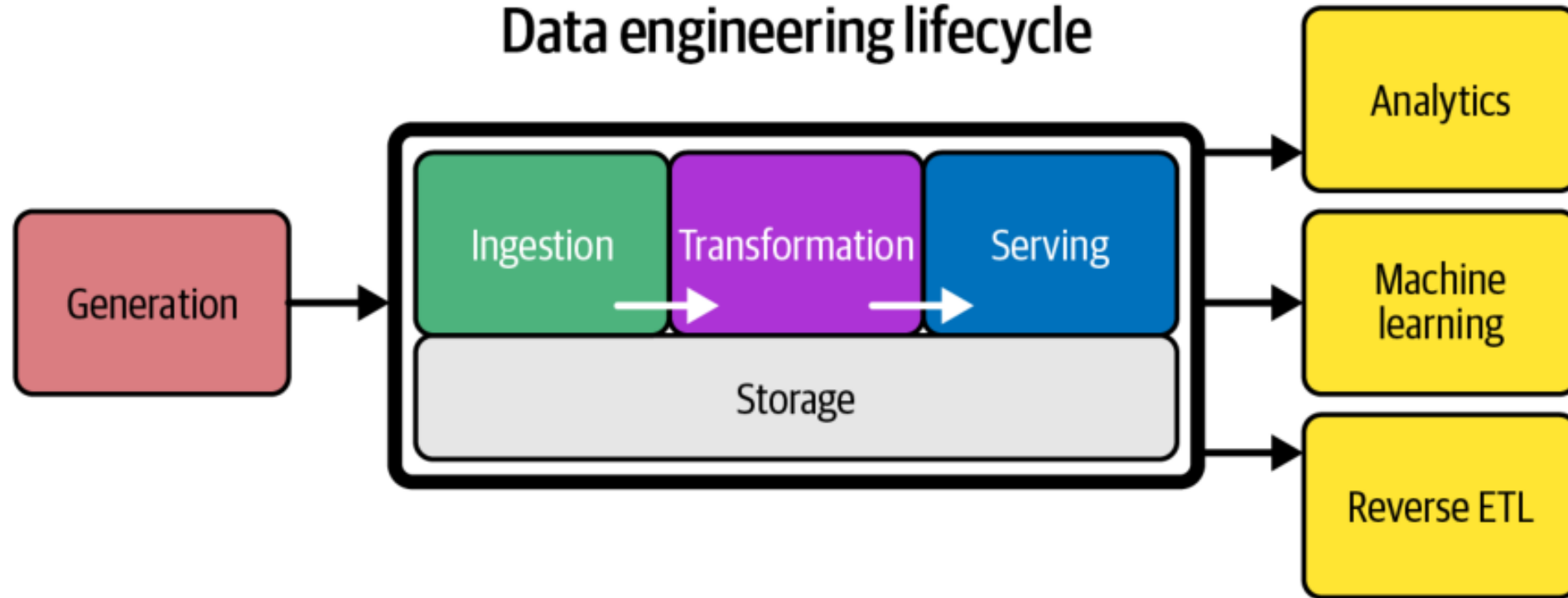
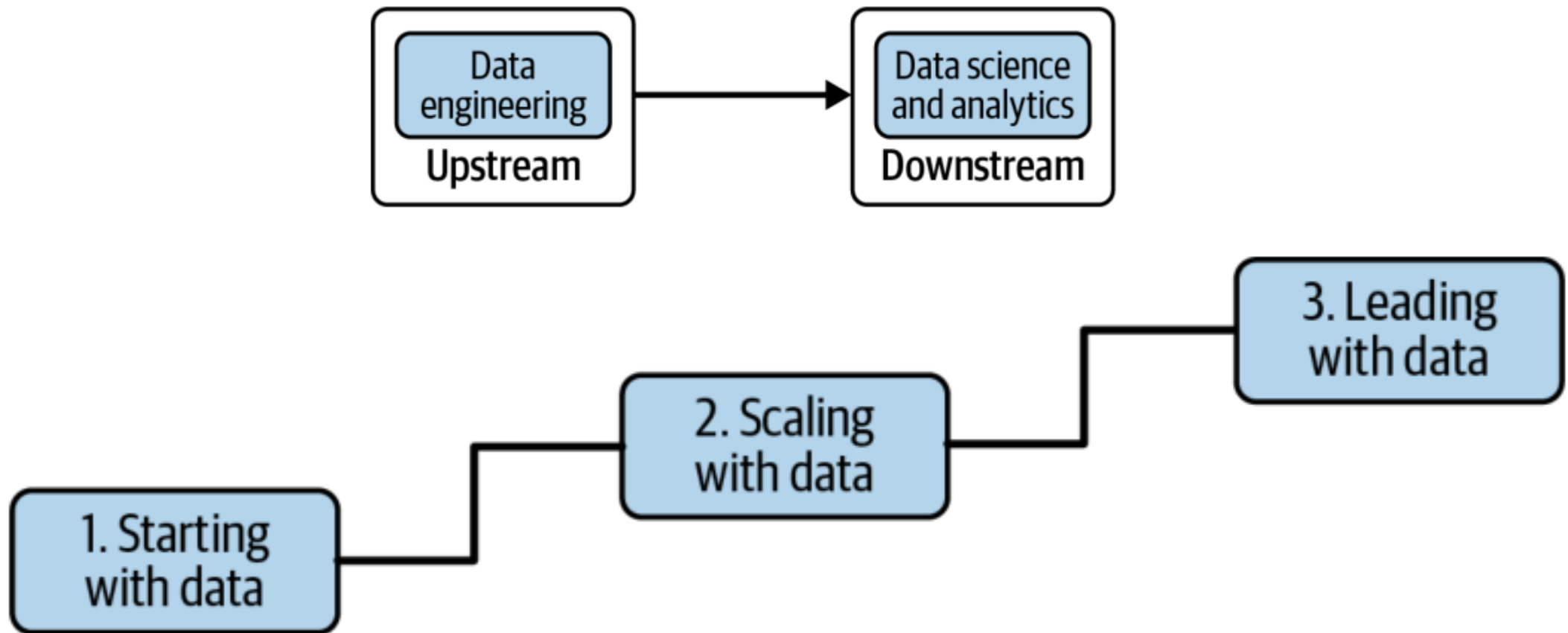


DATA ENGINEERING THEORY AND ADVICES

A HIGH LEVEL VIEW OF DATA ENGINEERING ECOSYSTEM:



UP AND DOWNSTREAM AND DATA MATURITY MODEL:



A DATA ENGINEER MUST UNDERSTAND:

1. داده: جنبه‌های مختلف داده و مدیریت آنرا
بشناسد.

2. تکنولوژی: ابزارهای مختلف کار با دیتا و
ارتباطشان با یکدیگر را درک نماید.



TWO TYPES OF DATA ENGINEERS:

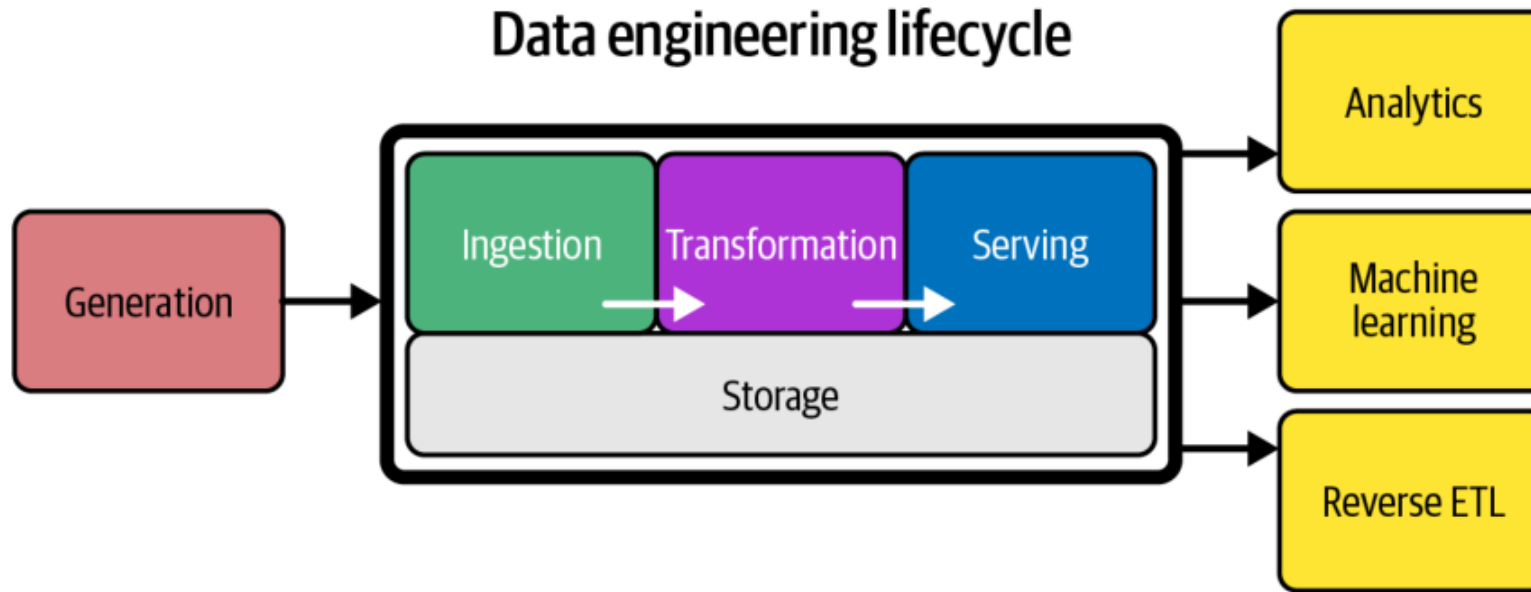
TYPE A (A STANDS FOR ABSTRACTION):

مهندسين داده در اين سطح و مرحله، از انجام کارهای سخت و پیچیده پرهیز می‌نمایند. آنها تلاش می‌کنند که معماری‌های مختلف داده را تا حد ممکن بصورت مجرد (روی کاغذ پیاده کنند که به مرحله‌ی واقعیت هنوز نرسیده است) تهیه و تدوین نمایند. این سطح از مهندسين عمومین چیزی را ابداع نمی‌کنند و از ابزارهای آماده استفاده می‌کنند.

TYPE B (B STANDS FOR BUILD):

سازندگان اصلی پایپ‌لاین‌ها، این سطح از مهندسين هستند. در اکثر موارد به ابداع سیستم‌های جدید می‌پردازند.

FIVE STAGES IN DATA ENGINEERING LIFECYCLE:

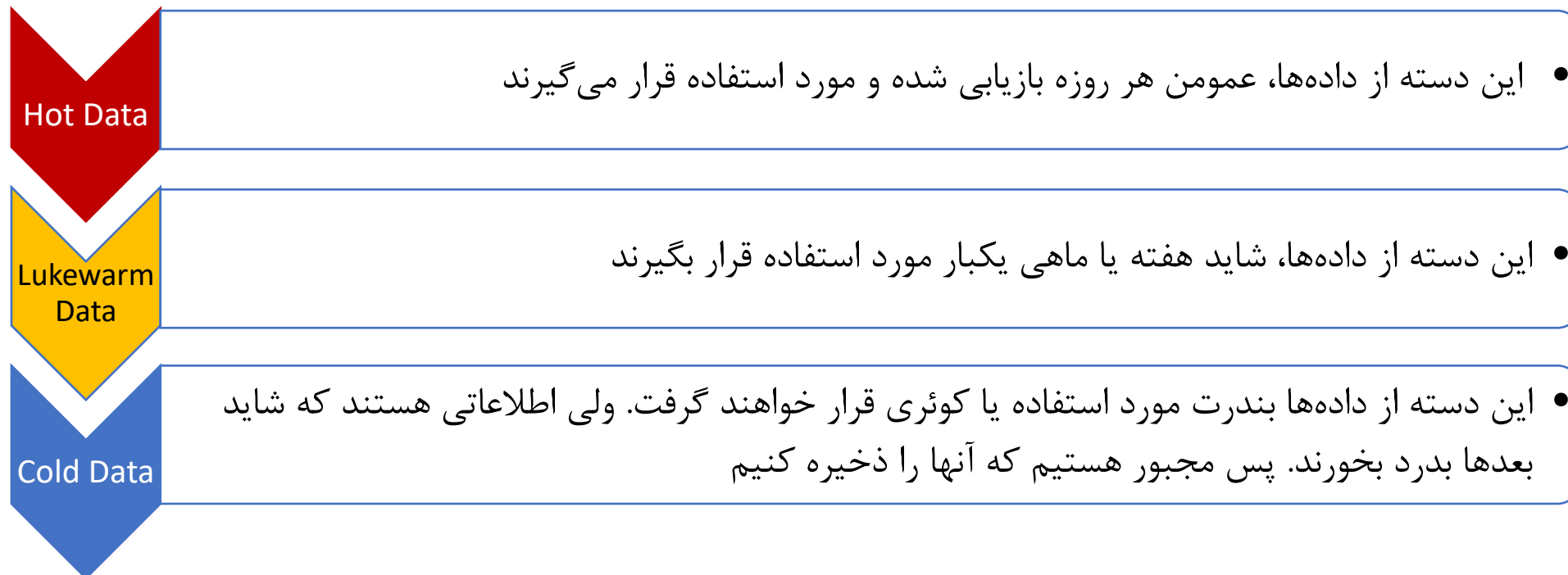


1- GENERATION: SOURCE SYSTEMS:

منبع داده می‌تواند یک دستگاه IoT، یک دیتابیس، کلیک‌های یک وبسایت و یا هر چیز دیگر باشد. مهندسین داده باید درک درستی از منابع داده پیدا کنند و نحوه‌ی کار با آنها را بخوبی بشناسند. دقت کنید که مهندس داده مالک یا کنترل کننده‌ی منابع دیتا نیست.

2- STORAGE:

یکی از پیچیده‌ترین (Most Complicated) مراحل در مهندسی داده، مرحله‌ی ذخیره‌سازی داده‌هاست. چون تا حد زیادی به پول و مصرف منابع مالی مرتبط است. در فاز ذخیره‌سازی داده، باید فراوانی دسترسی به دیتا (Data Access Frequency) را خوب بفهمیم.



DATA STORAGE IN CLOUD ENVIRONMENT IS VERY CHEAP
BUT DATA RETRIEVAL IS EXPENSIVE!

انتخاب راه حل برای ذخیره سازی خیلی مهم است. این مورد به استفاده های ما از دیتا، حجم داده، فراوانی دریافت (هضم) داده و فرمت دیتاهای ورودی وابسته است.

3- INGESTION:

بیشترین گیر و خطا (یا حتی بن بست) در مرحله ی هضم دیتا از منابع مختلف داده رخ می دهد.

بطور کلی در مرحله‌ی هضم، دو مفهوم مهم وجود دارد

BATCH:

بعنوان مثال در پایان روز کاری، اطلاعات فروش را در آن روز مورد پردازش قرار می‌دهیم (هضم بصورت بچ از منابع دیتا (دیتابیس) به سیستم پردازشی).

STREAMING:

انتقال دیتا به سیستم‌های پائین‌دستی (Downstream) بصورت پیوسته و در لحظه (دیتای در لحظه، یعنی بلافاصله پس از تولید به پائین‌دستی‌ها برسد)

زمانی که منابع محاسباتی خوب نبودند، پردازش بصورت بچ فراگیر بود. اما امروزه دریافت و پردازش استریم طرفداران بیشتری دارد.

در بخش هضم، دو واژه‌ی دیگر نیز وجود دارد و حائز اهمیت می‌باشند:

PUSH VS. PULL:

در مدل پوش، یک منبع داده، دیتا را بسمت هدف ارسال می‌کند (می‌نویسد). هدف می‌تواند یک دیتابیس، دیتالیک و یا یک سیستم فایلی باشد. در مدل پول، داده از منبع دیتا بازیافت (retrieved) می‌شود.

4- TRANSFORMATION:

در فاز تبدیل، داده‌های ما به ارزش برای مصرف کننده‌ی پائین‌دستی بدل می‌گردند. چون پائین‌دستی‌ها از هر داده‌ای نمی‌توانند بهره ببرند. در یک سازمان، به محض اینکه دیتاساینטיست‌ها تشخیص دادند که از چه ویژگی‌هایی از دیتا، قرار است استفاده نمایند، تیم مهندسی داده فرایند استخراج این ویژگی‌ها و در اختیار قرار دادن آنها به سیستم تحلیل را بصورت اتوماتیک تبدیل می‌نمایند.

5- SERVING DATA

در اختیار مصرف‌کنندگان قرار دادن دیتا (سرو دیتا) یکی از هیجان‌انگیزترین فازها در مهندسی داده است. باید در این فاز بررسی کنیم که چه کسانی و چه چیزهایی ممکن است مصرف‌کننده‌ی داده باشند. چند مورد را در ادامه می‌بینیم:

I. ANALYTICS

II. ML

III. REVERSE ETL (NEW TECH AND TREND)

یک مثال از ای‌تی‌ال وارونه اینگونه است؛ دیتا از سمت دیتاورهوز بسمت تیم مارکتینگ می‌رود برای استفاده‌ی بهینه، در واقع برای ارسال ایمیل‌های تبلیغاتی برای مشتریان متعهد سازمان.

DESIGNING A GOOD DATA ARCHITECTURE

معماری داده، شرحی از ساختارها و تعاملات مولفه‌های اصلی و منابع دیتایی یک سازمان است.

“NEVER SHOOT FOR THE BEST ARCHITECTURE, BUT RATHER THE LEAST WORST
ARCHITECTURE”

MARK RICHARDS AND NEAL FORD, FUNDAMENTALS OF SOFTWARE ARCHITECTURE

یک معماری داده‌ی خوب، در خدمت بیزینس و سازمان است و مولفه‌هایی دارد که قابل استفاده‌ی مجدد هستند. در حالی که یک تعادل مطلوبی بین عملکرد و هزینه را برقرار می‌سازد.

9 PRINCIPLES OF DATA ENGINEERING ARCHITECTURE

1- CHOOSE COMMON COMPONENTS WISELY:

نیاز به ابداع و طراحی مولفه‌های جدید نیست. عمومن تیم‌ها تشویق می‌شوند که از مولفه‌های معمول و مرسوم استفاده کنند.

2- PLAN FOR FAILURE

حتمن بحث failure یا مشکل را در طراحی خود در نظر بگیرید. بقول آقای ووگل (مدیر تکنولوژی آمازون):

EVERYTHING FAILS, ALL THE TIME

موارد زیر را در طراحی خود (در این فاز) در نظر بگیرید:

AVAILABILITY, RELIABILITY, RECOVERY TIME OBJECTIVE(RTO), RECOVERY POINT OBJECTIVE (RPO)

RTO: بیشترین زمان قابل قبول برای از دسترس خارج شدن یک سرویس

RPO: بیشترین مقدار از دسترفت داده بعد از هر اتفاق بد

3- ARCHITECT FOR SCALABILITY:

بار فعلی سیستم‌های خود را اندازه‌گیری کنید. نوسانات و افزایش‌های شدید را تخمین بزنید و با توجه به این موارد بار وارده را در چند سال آتی برآورد کنید. اینگونه بهتر می‌توانید یک سیستم و معماری طراحی کنید که پایدارتر و کاراتر و مقیاس‌پذیرتر باشد.

4- ARCHITECT IS LEADERSHIP:

5- ALWAYS BE ARCHITECTING:

6- BUILD LOOSELY COUPLED SYSTEMS:

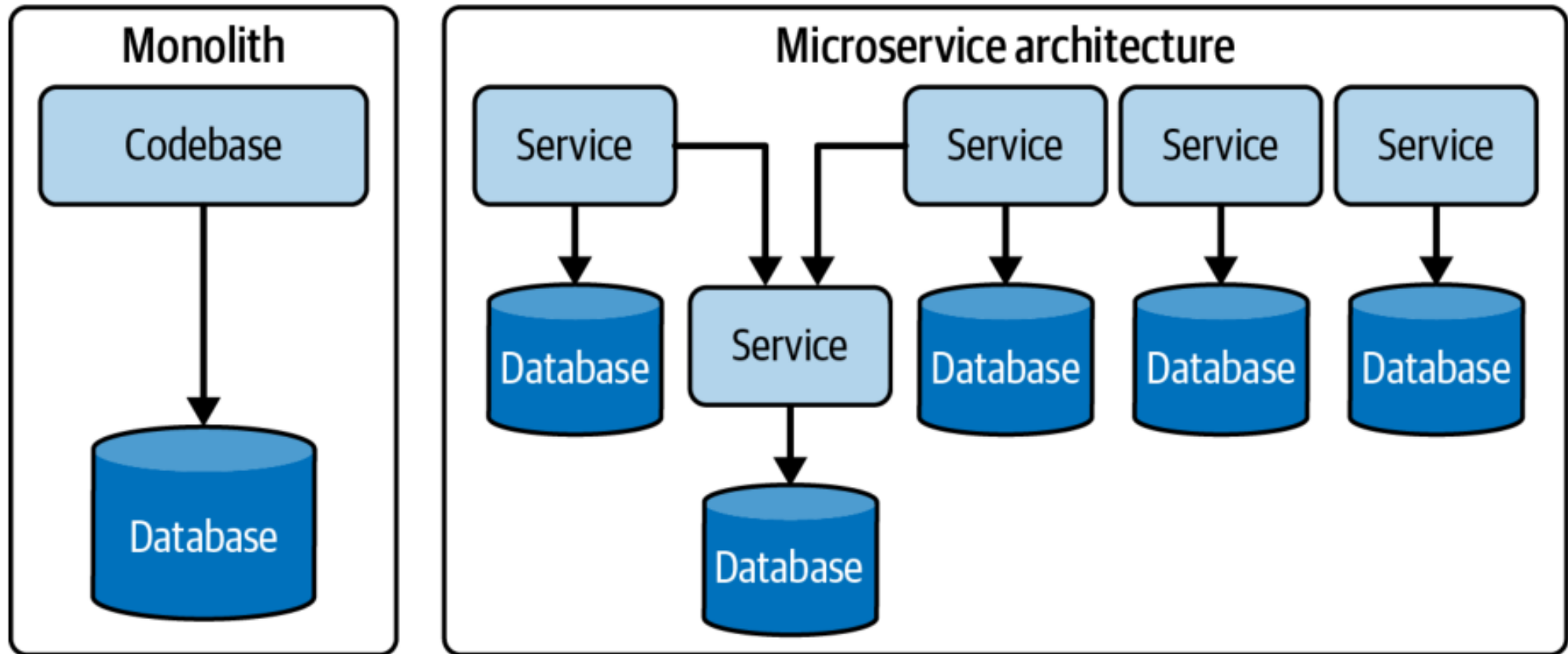
7- MAKE REVERSIBLE DECISIONS:

8- PRIORITIZE SECURITY:

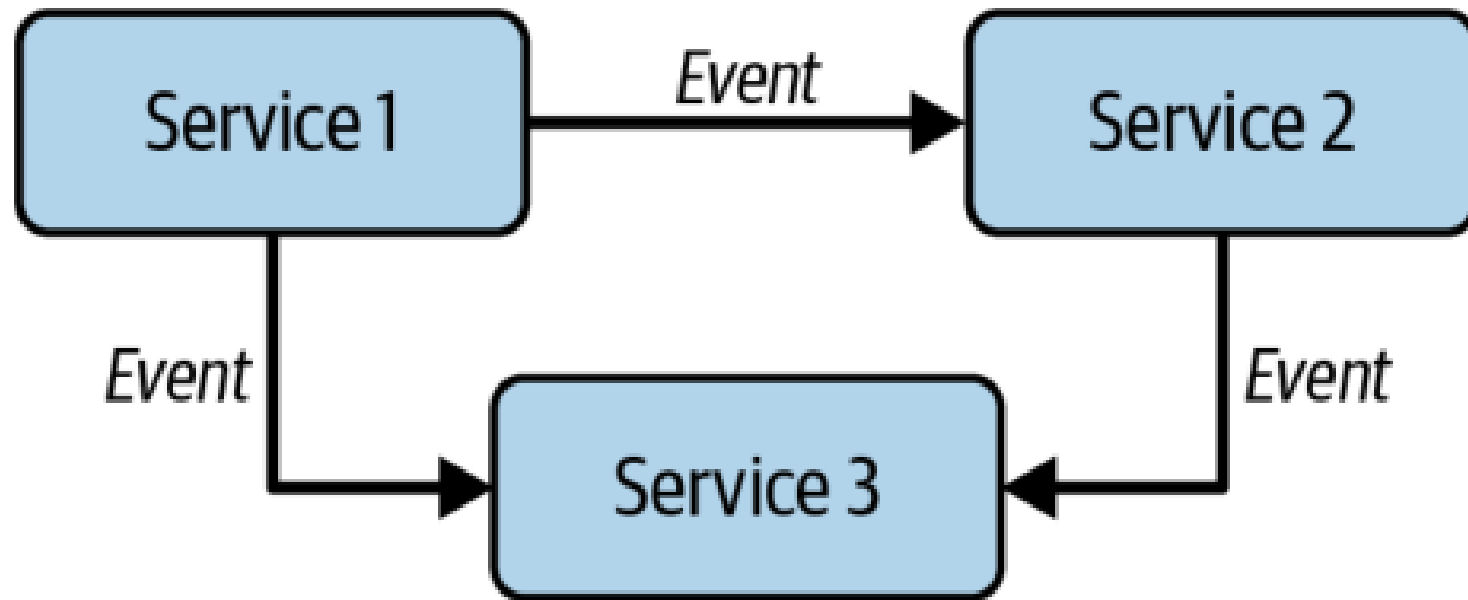
9- EMBRACE FINOPS:

خرج و مخارج کلاذ را درست مدیریت کنید.

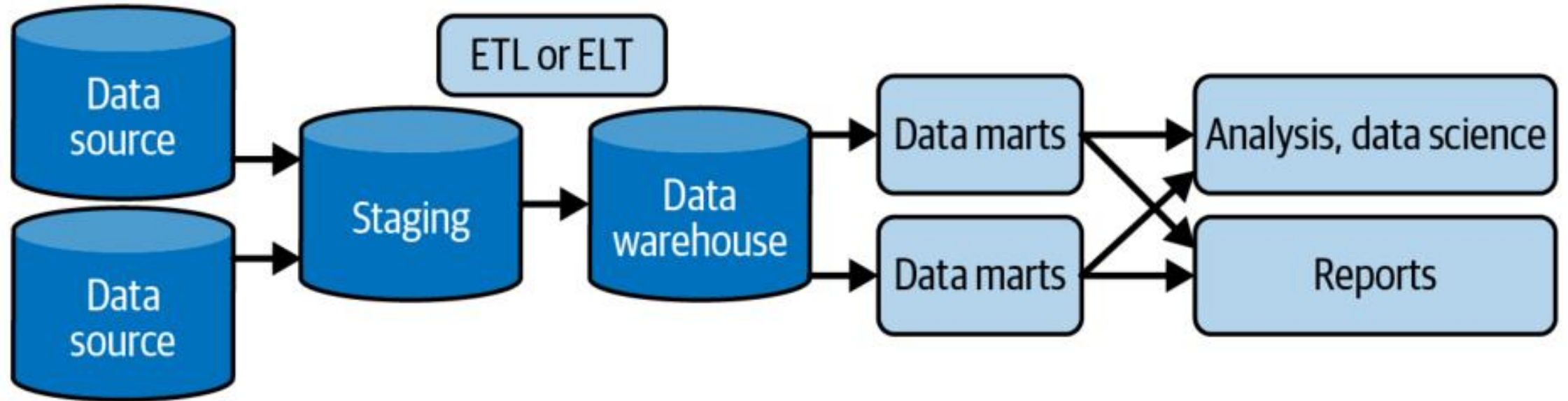
TWO MAIN ARCHITECTURES:



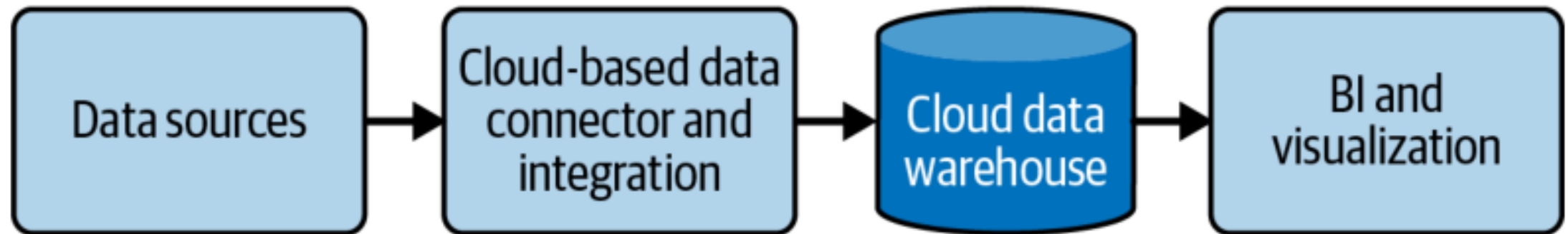
EVENT-DRIVEN ARCHITECTURE:



DATA WAREHOUSE AND DATA MARTS:



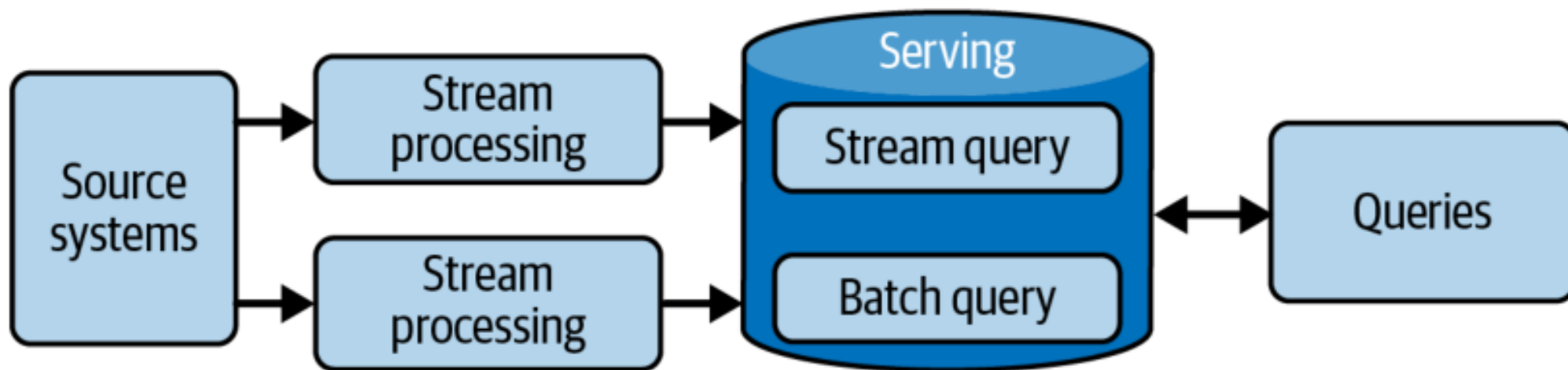
MODERN DATA STACK:



LAMBDA ARCHITECTURE:

در معماری لَمدا، ما هر سه مولفه‌ی بچ، استریم و سرو داده را در کنار هم داریم که بطور مستقل از یکدیگر، کار می‌کنند.

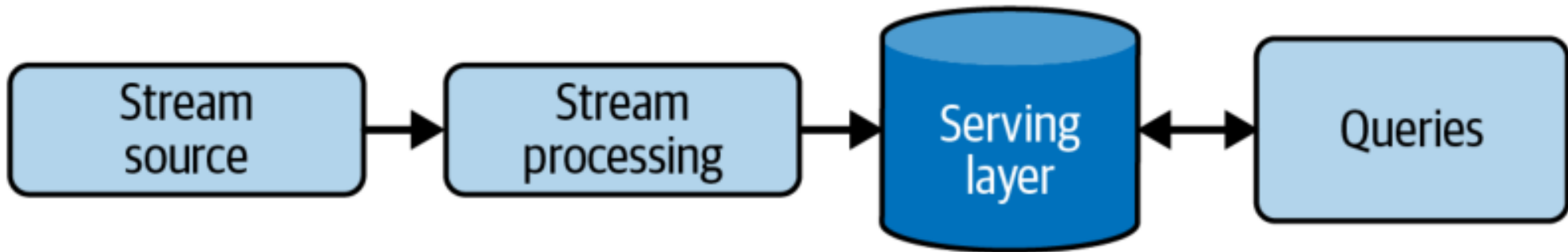
عیب: مدیریت دو سیستم بچ و استریم همزمان سخت است



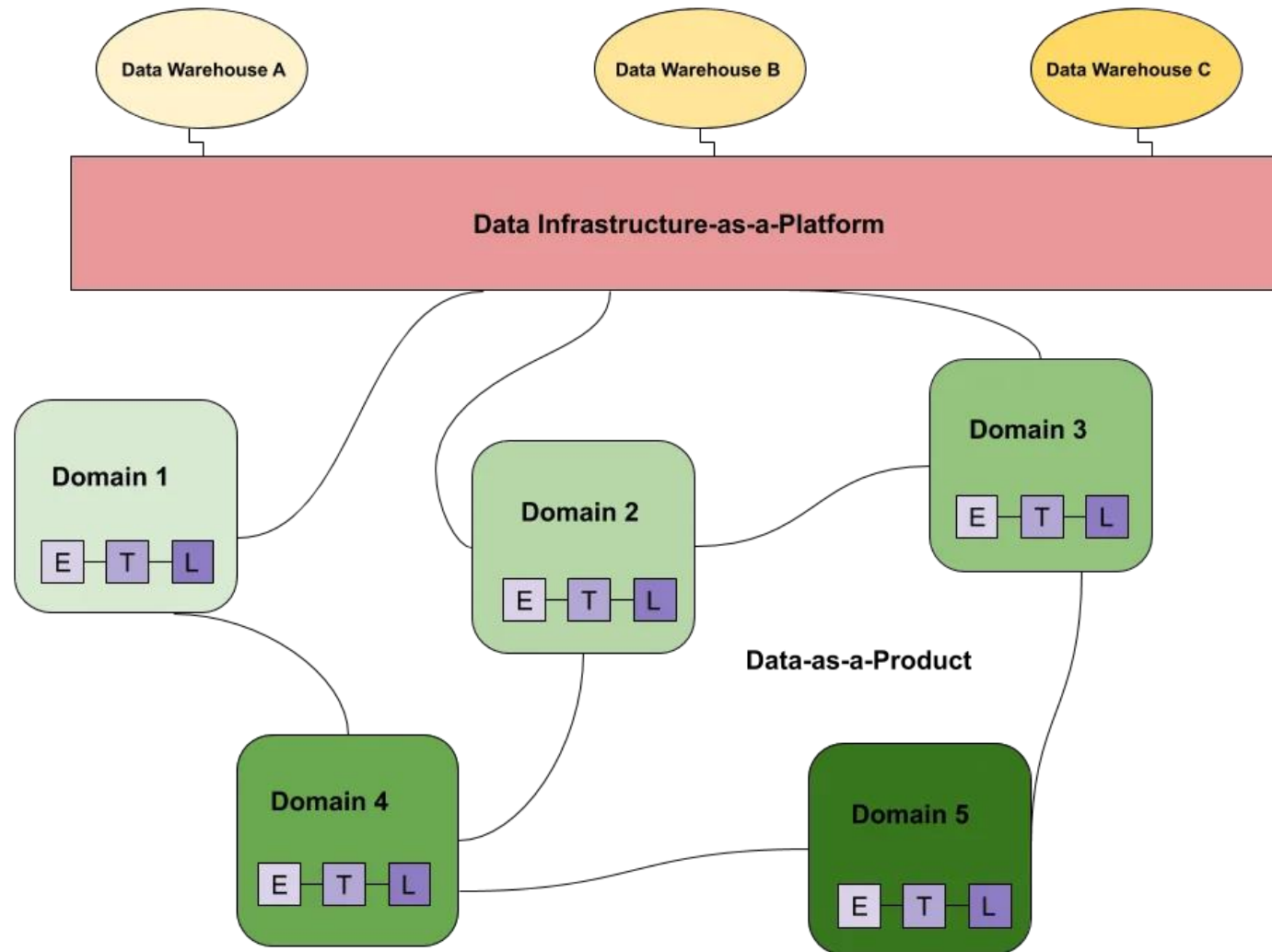
KAPPA ARCHITECTURE:

در معماری کاپا، ما تنها از یک پلتفرم پردازش استریم برای هندل کردن تمام کارها استفاده می‌کنیم.

عیب: سخت بودن در عمل و گران بودن در پیاده‌سازی



DATA MESH ARCHITECTURE:



DIFFERENCE BETWEEN ARCHITECTURE AND TOOL:

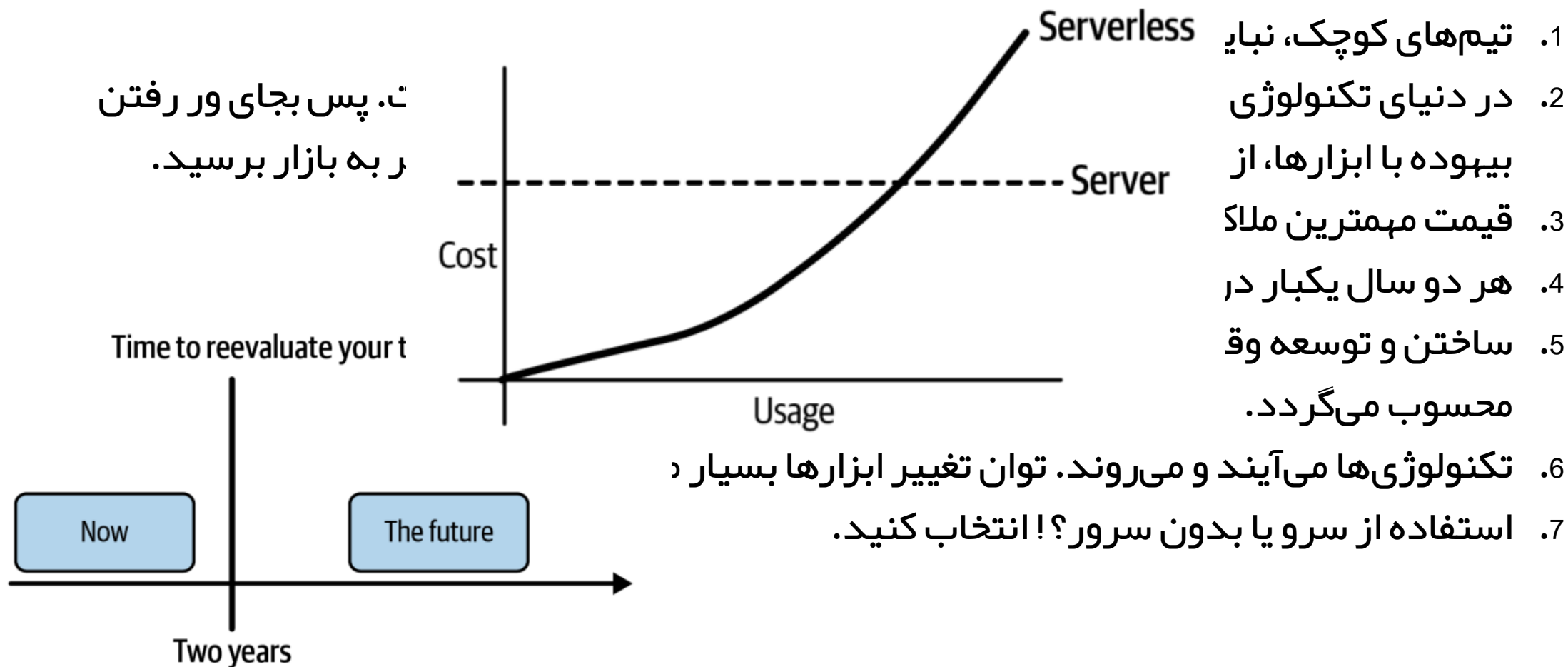
Architecture

ARCHITECTURE IS THE WHAT, WHY AND WHEN.

Tools

TOOLS ARE USED TO MAKE THE ARCHITECTURE A REALITY. TOOLS ARE THE HOW.

SOME PIECES OF ADVICE:



SOME PIECES OF ADVICE:

1. تیم‌های کوچک، نباید کارهای پیچیده‌ی کمپانی‌های بزرگ را تقلید کنند.
2. در دنیای تکنولوژی سرعت ورود به بازار (Speed to market) همیشه برنده است. پس بجای ور رفتن بیهوده با ابزارها، از همان تکنولوژی‌های روز و مرسوم استفاده کنید تا زودتر به بازار برسید.
3. قیمت مهمترین ملاک برای انتخاب تکنولوژی است.
4. هر دو سال یکبار در استفاده از ابزارهای خود تجدید نظر کنید.
5. ساختن و توسعه وقتی معنی دارد که برای کسب و کار سودمند باشد. در غیراینصورت اتلاف وقت محسوب می‌گردد.
6. تکنولوژی‌ها می‌آیند و می‌روند. توان تغییر ابزارها بسیار مهم هستند.
7. استفاده از سرو یا بدون سرور؟! انتخاب کنید.

(FOUR + ONE) PLACES TO RUN OUR TECHNOLOGIES:

1- ON PREMISES

2- CLOUD

3- HYBRID CLOUD (ON PREM + CLOUD)

4- MULTICLOUD (AWS + AZURE + GCP +...)

DISADVANTAGE: HANDLING DIFFERENT PANEL IS HARD!

FIVE- CLOUD OF CLOUDS (SNOWFLAKE- IT OFFERS JUST ONE SINGLE
PANEL)

KNOW THREE AWESOME GUYS IN THE WORLD OF DE AND CLOUD

1- JOE REIS

2- STEPHEN MAARAK (+2M STUDENTS IN UDEMY= 100 BILION TOMAN 🤖)



Mahdie Panahian • 1st

Data Analyst

...

من مدتی قبل گزارش دیجی‌پی رو توی صفحه‌ام منتشر کردم و از بقیه خواستم ایرادات فنی این گزارش رو بگن. الان با این پست مواجه شدم. اول به دیجی‌پی تبریک میگم بابت این حجم از پیگیری و انتقاد پذیریشون و دوم به خودم، چون احساس میکنم به عنوان یک دیتا انالیست تازه نفس، تا حدودی به چیزی که همیشه میخواستم رسیدم اون هم بی اعتنا بودن و گذر نکردن از جزییات هر چند کوچیکه. در نهایت بازم آرزوی موفقیت بیشتر دارم برای این سازمان و کارمندانش.

1. تخصصی کار کنید

2. کلاس را جدی بگیرید

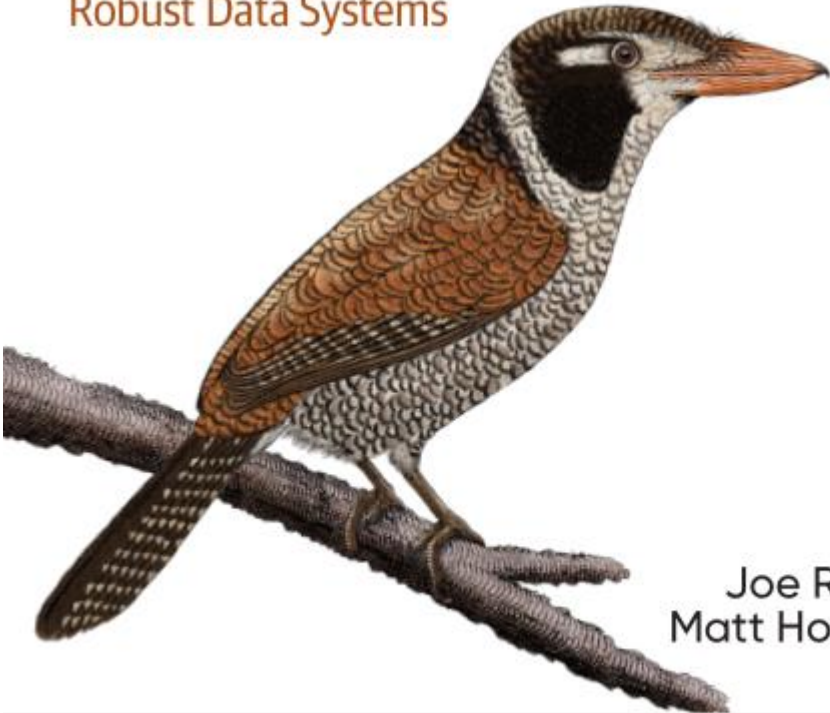
3. جاوا یاد بگیرید

4. از کار بزرگان ایراد بگیرید و نظر دهید

O'REILLY®

Fundamentals of Data Engineering

Plan and Build
Robust Data Systems



Joe Reis &
Matt Housley

THIS TALK IS BROUGHT TO YOU
USING THE FIRST FOUR
CHAPTERS (OUT OF ELEVEN) OF
THIS AWESOME BOOK.