

Big Data Course

Chap 1- Introduction

Electrical Engineering department of
Amirkabir University of technology

Dr. Mohammadreza Pourfard

August 2025

صفحات درس و استاد



اتاق: طبقه سوم ساختمان ابوریحان، دانشکده مهندسی برق، آزمایشگاه سیستمهای هوشمند دیجیتال

@Electronic_daneshbonyan



pourfardm@gmail.com



+982164543373



Research Groups

In today's rapidly evolving technological landscape, numerous research groups are at the forefront of Innovation, each focusing on transformative areas that have the potential to reshape Industries and enhance our daily lives. Among these, the Internet of Things (IoT) research group is dedicated to exploring the interconnectedness of devices and sensors, striving to create smart environments that foster efficiency and convenience. Meanwhile, our Artificial Intelligence (AI) team is delving into advanced algorithms and machine learning techniques to develop systems that can learn, adapt, and make informed decisions, pushing the boundaries of automation and intelligence. Additionally, the Big Data research group is



Laboratory page



Profile Contents

Home

Courses

- About
- Theses
- Research Interests
- Research Groups
- Research Projects
- Current Students
- Employment Records
- Contact
- Publications
- News

About

Adjunct Professor of Amirkabir University of Technology

Theses

B.S Degree:

Simulating Multi-Level Cache Memory

M.S Degree:

Detection and Tracking of a Human in different positions of its body

PhD Degree:

Texture Analysis and Separation for Characterization of Material's Structure through

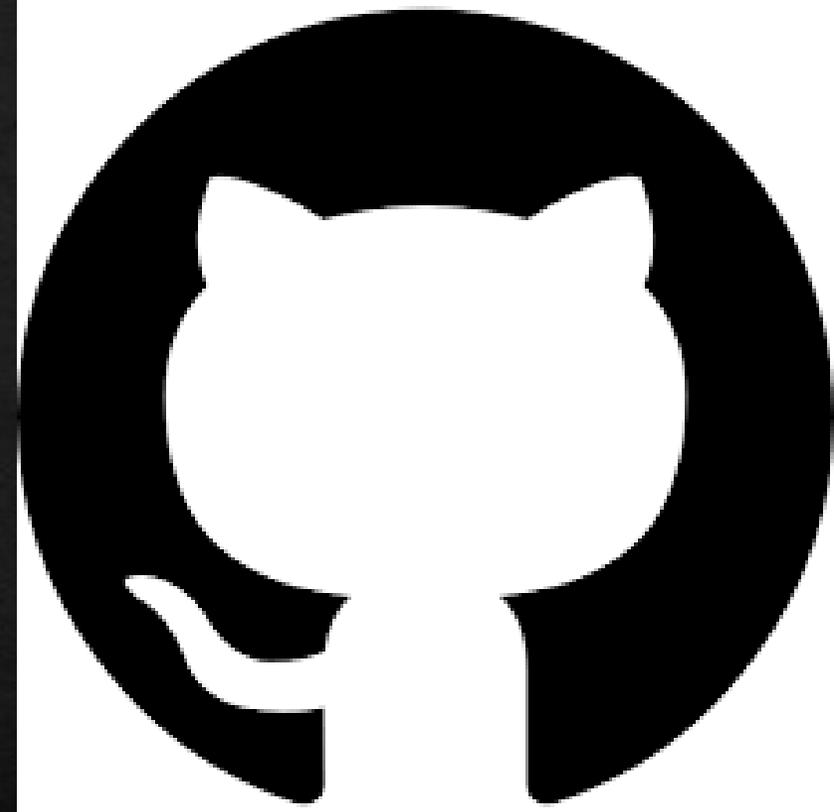
Research Interests

- Genetics data processing,

کانالها و گروه‌های مرتبط با درس

ADVANCED CODING

- AI -



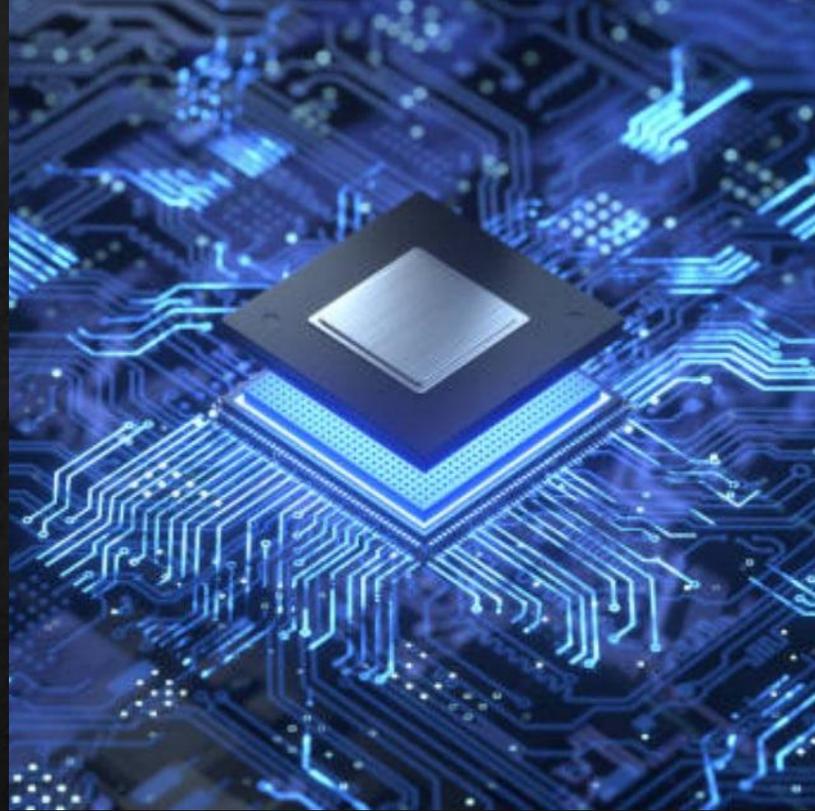
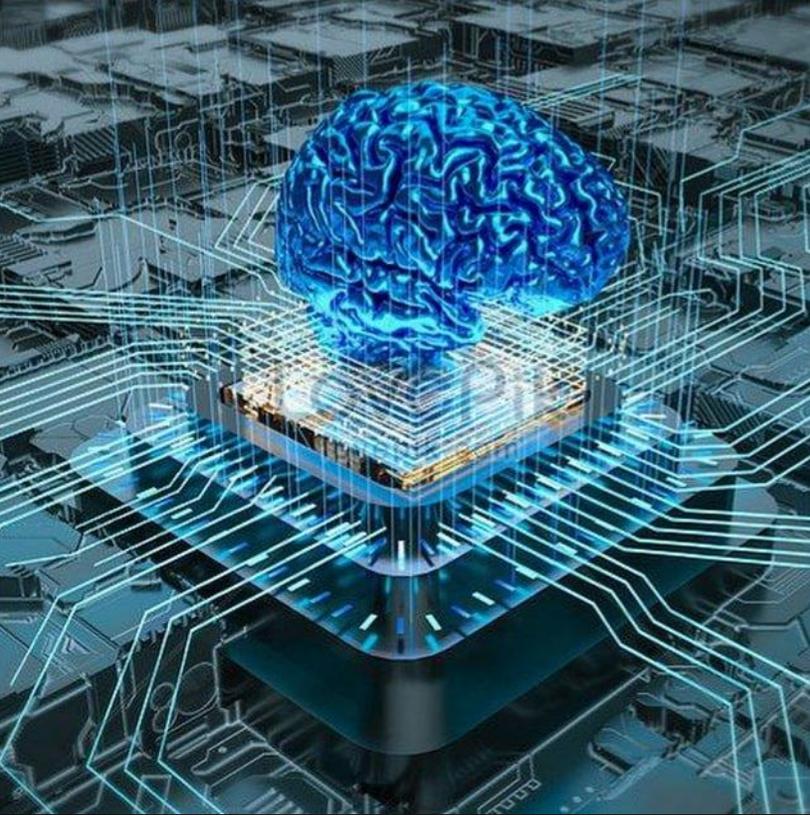
Telegram Channel:
https://t.me/Advanced_programming_algorithm_c



Telegram Group:
https://t.me/Advanced_Programming_Algorithm



GitHub Page:
<https://github.com/Pourfardm/>



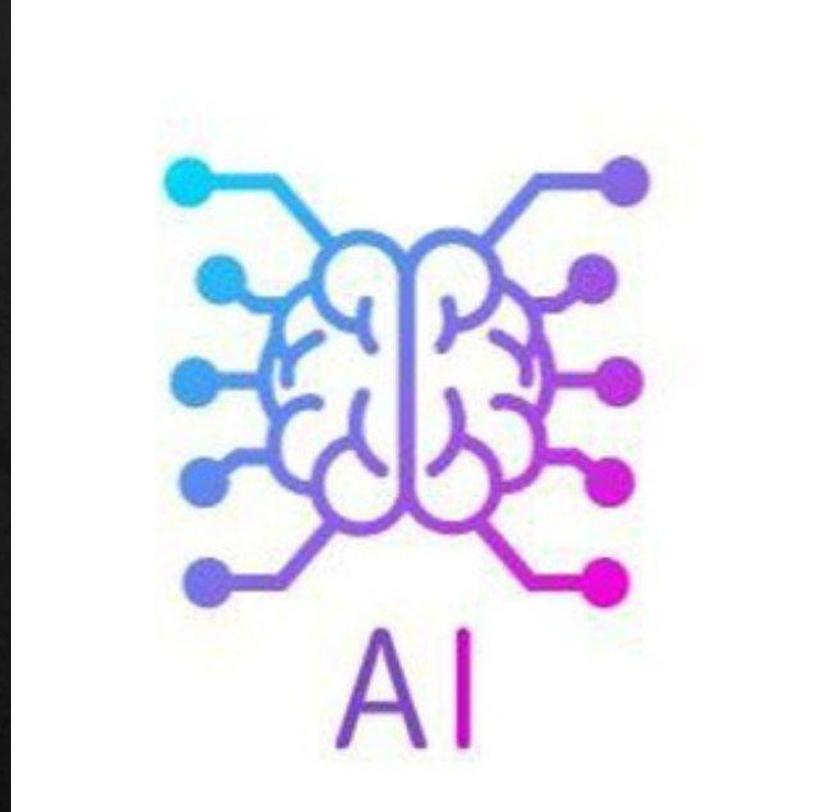
Telegram Channel:
[https://t.me/fpga_digital_logic_d
esign](https://t.me/fpga_digital_logic_design)



Telegram Group:
[https://t.me/Advanced_digital_a
nalog_design](https://t.me/Advanced_digital_analog_design)



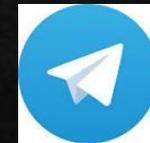
YouTube Channel:
[https://www.youtube.com/chan
nel/UCa8ZFUxp37Vy-
KEbx7SbbHQ](https://www.youtube.com/channel/UCa8ZFUxp37Vy-KEbx7SbbHQ)



Telegram group:
https://t.me/AI_Bigdata_amirkabir_group



Telegram Channel:
https://t.me/Deeplearning_BigData_AI



Private Telegram Group:
For class student



Telegram Channel:

https://t.me/AI_Neuroscience_Genomic_Data

YouTube Channel:

<https://www.youtube.com/channel/UCa8ZFUxp37Vy-KEbx7SbbHQ>

اهداف درس،

اهداف کلیدی درس

- ❖ آشنایی با کاربردهای مفید پردازش کلان داده به منظور راه اندازی کسب و کار
- ❖ آشنایی با ابزارهای مورد استفاده در شرکتهای دانش بنیان برای پردازش کلان داده
- ❖ آشنایی با الگوریتم های جدید و پیشرفته پردازش کلان داده
- ❖ برقراری ارتباط بین مهارت آشنایی با سخت افزار همانند GPU

مروی بر نسل دانشگاهها در جهان
اهمیت وروده دانشگاه نسل سوم و
ارایه درس متناسب با دانشگاه نسل سوم

مقایسه نسل های مختلف دانشگاهها در جهان

نسل دانشگاه	دوره تاریخی	مأموریت اصلی	ویژگی های کلیدی	نمونه ها	نقش اجتماعی
نسل اول – آموزش محور	قرون وسطی تا قرن ۱۹	آموزش و انتقال دانش	- تمرکز بر علوم کلاسیک (فقه، فلسفه، پزشکی) - نبود پژوهش سیستماتیک - متون سنتی و اقتدار دینی	بولونیا، پاریس، آکسفورد	تربیت نخبگان و مشروعیت بخشی به دین و حکومت
نسل دوم – پژوهش محور	قرن ۱۹ تا اواسط قرن ۲۰	آموزش + پژوهش	- پژوهش بنیادی - آزادی آکادمیک - شکل گیری رشته های علمی مدرن	برلین (هومبولت)، کمبریج، MIT (مدرن اولیه)	تولید علم و فناوری، پیشرفت علمی جوامع
نسل سوم – کارآفرین	اواخر قرن ۲۰	آموزش + پژوهش + نوآوری / کارآفرینی	- تجاری سازی دانش - مراکز رشد و پارک فناوری - ارتباط نزدیک با صنعت	استنفورد، MIT، کمبریج (Silicon Fen)	موتور اقتصاد دانش بنیان
نسل چهارم – اجتماعی / مسئولیت پذیر	اوایل قرن ۲۱	آموزش + پژوهش + کارآفرینی + مسئولیت اجتماعی	- توسعه پایدار - عدالت اجتماعی - توجه به محیط زیست و حکمرانی داده	دانشگاه های اسکاتلندی، تورنتو، برخی دانشگاه های آسیایی	حل مسائل کلان بشری (اقلیم، نابرابری، سلامت)
نسل پنجم – هوشمند/شبکه ای	دهه های اخیر (در حال شکل گیری)	یکپارچه سازی آموزش + پژوهش + نوآوری با فناوری دیجیتال	- هوش مصنوعی، big data، آموزش هوشمند - دانشگاه شبکه ای جهانی - یادگیری شخصی سازی شده	پروژه های دانشگاه های مجازی جهانی، edX، دانشگاه های آسیای شرقی	هدایت جامعه به سوی اقتصاد دیجیتال و تمدن AI



نکته مهم این درس

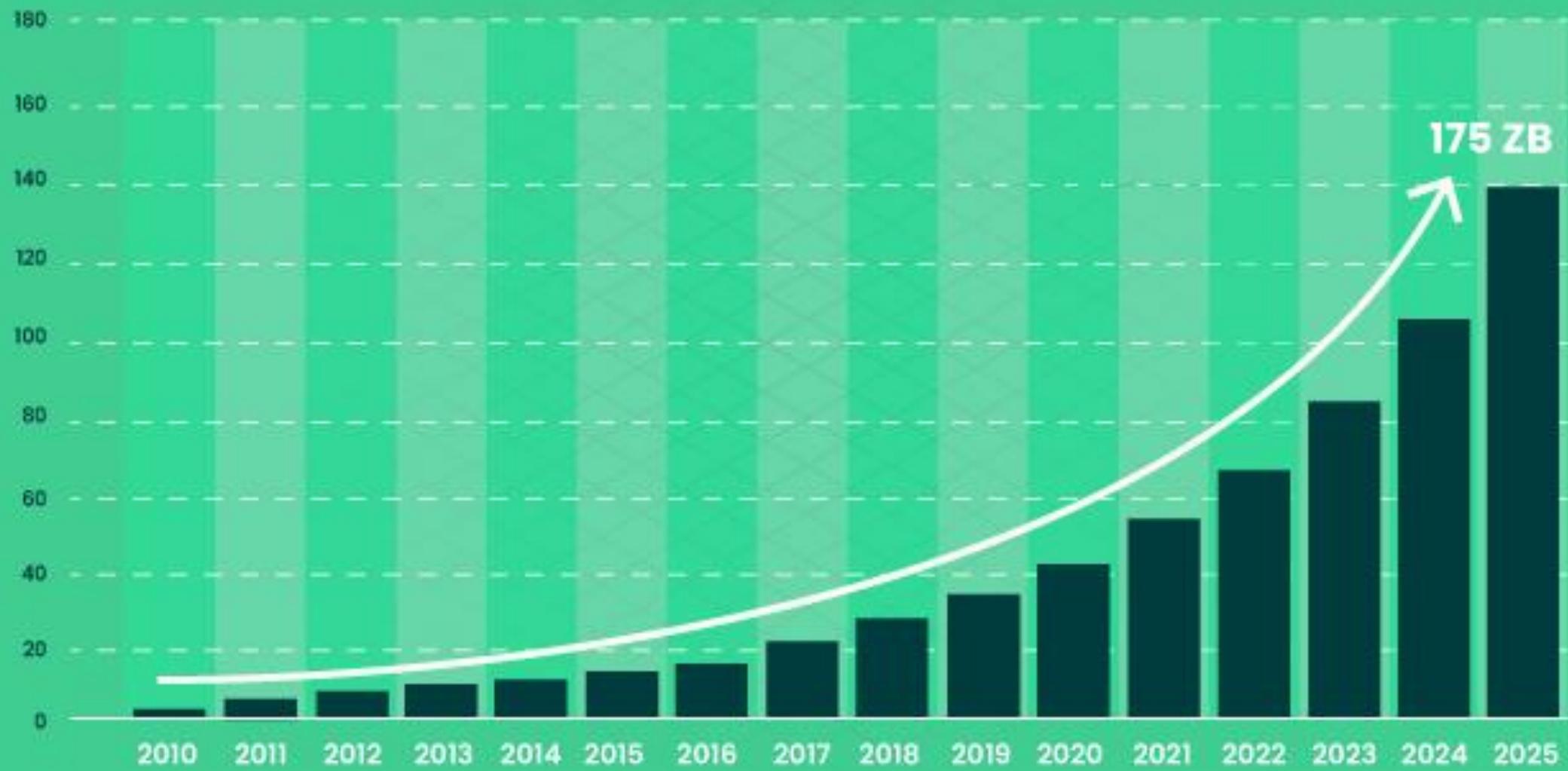
به دلیل رویکرد دانشگاه نسل سوم

مركز این درس فقط روی الگوریتم نیست

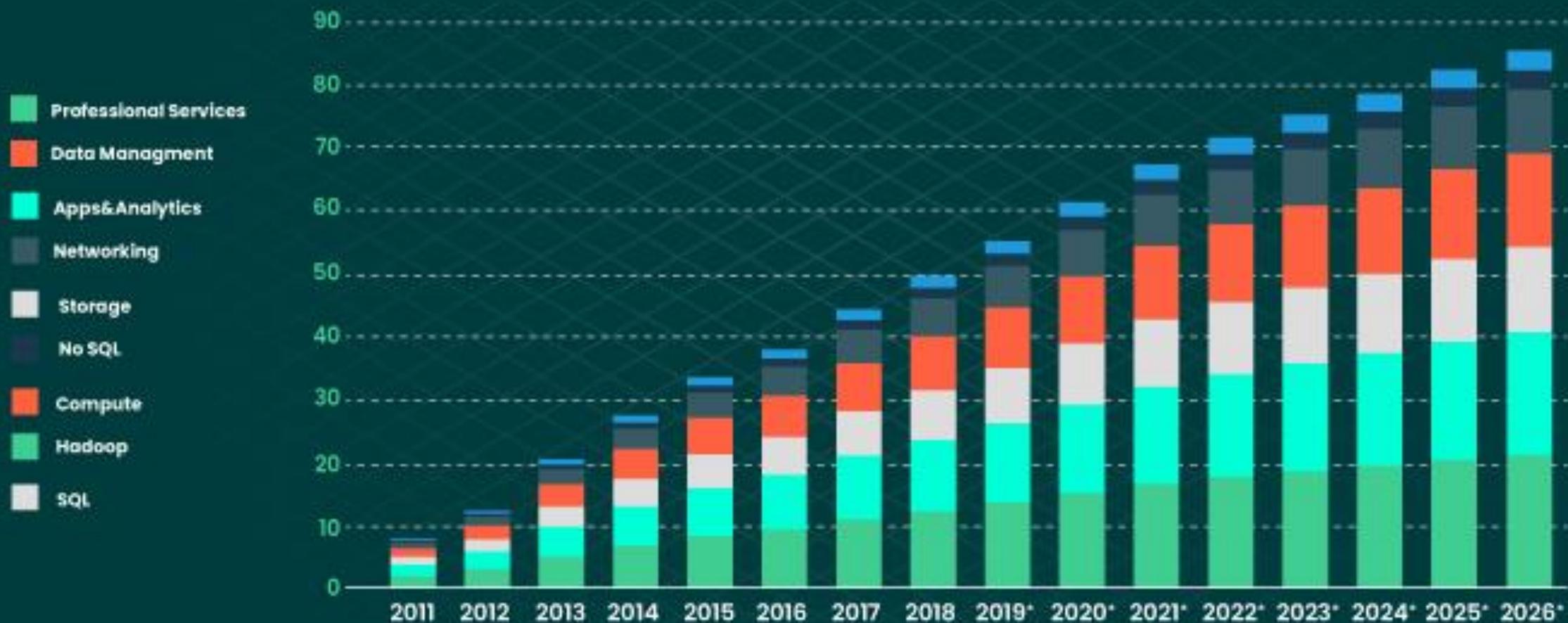
حجم و حسناك و خيره كننده توليد داده در همان

THE INTERNET IN 2023 EVERY MINUTE





Big Data Market Forecast Worldwide from 2011 to 2026, by segment (in billion U.S. dollars)



اہمیت BigData



Importance of Big Data for Business



کاربرد های BigData

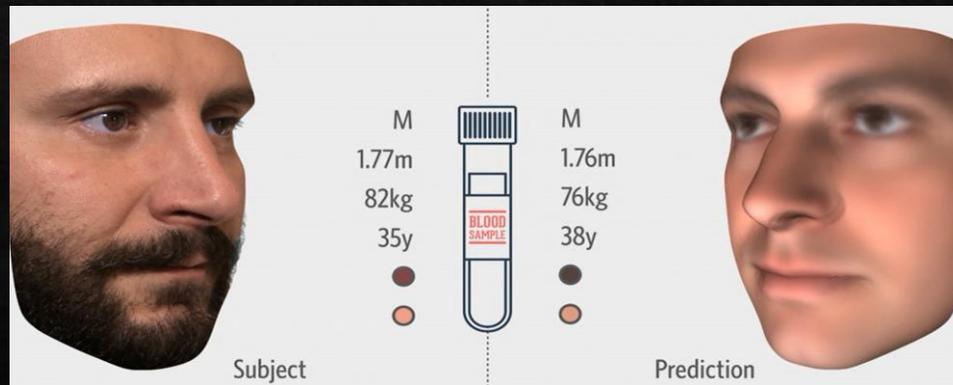
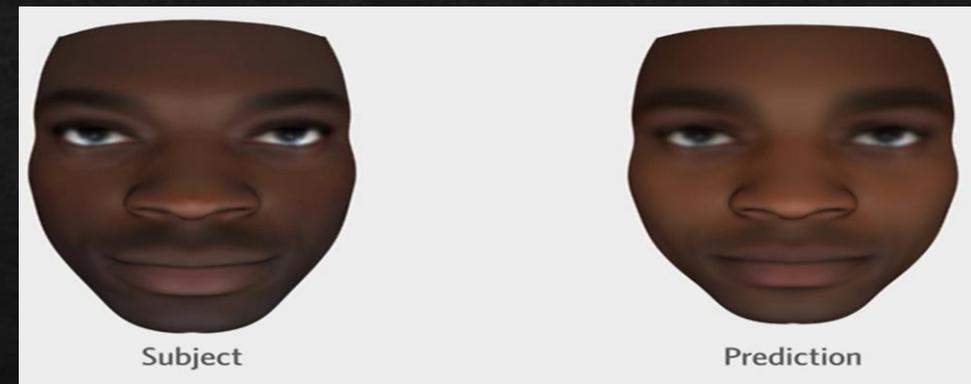
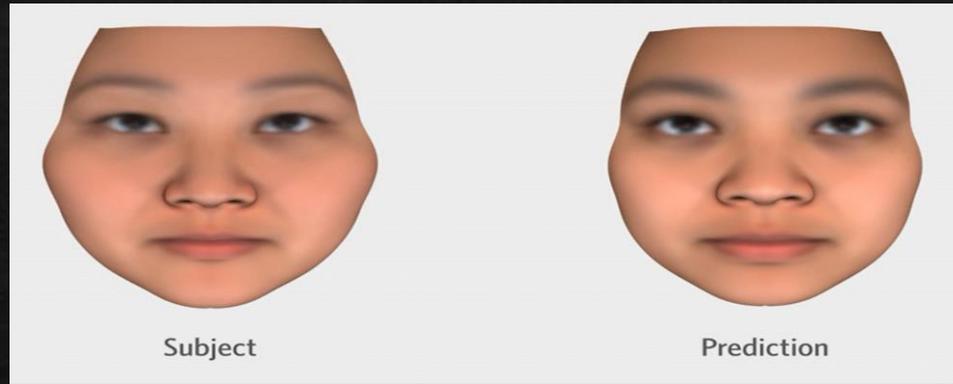
BIG DATA APPLICATIONS



Big Data Contributions to Healthcare

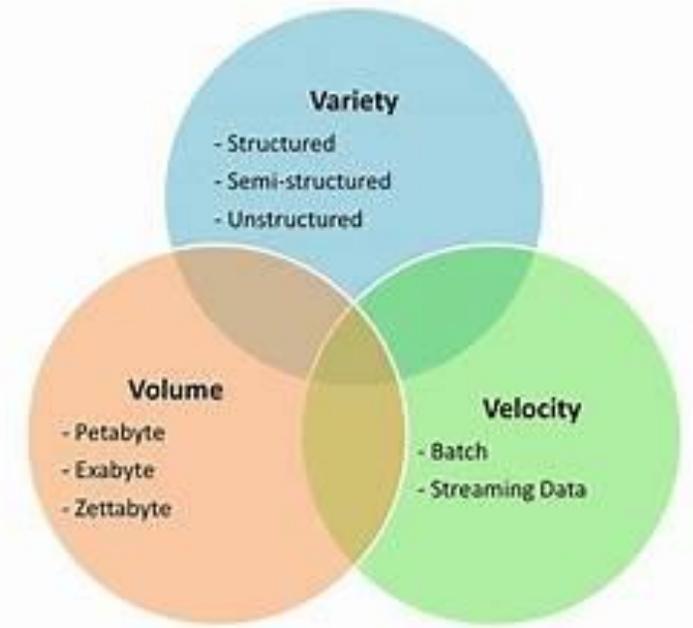
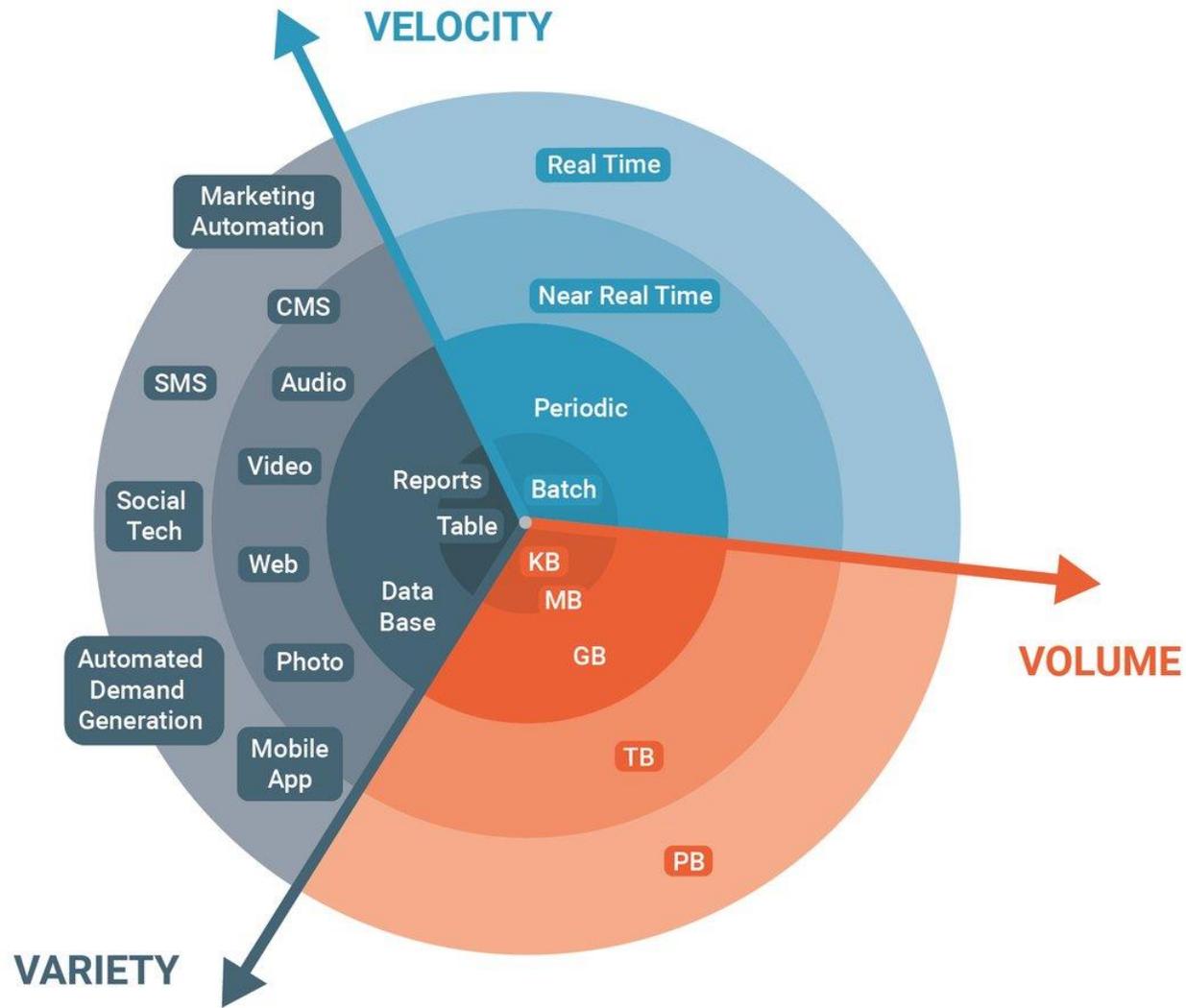


شرح مساله - پیش بینی چهره در سال ۲۰۱۶

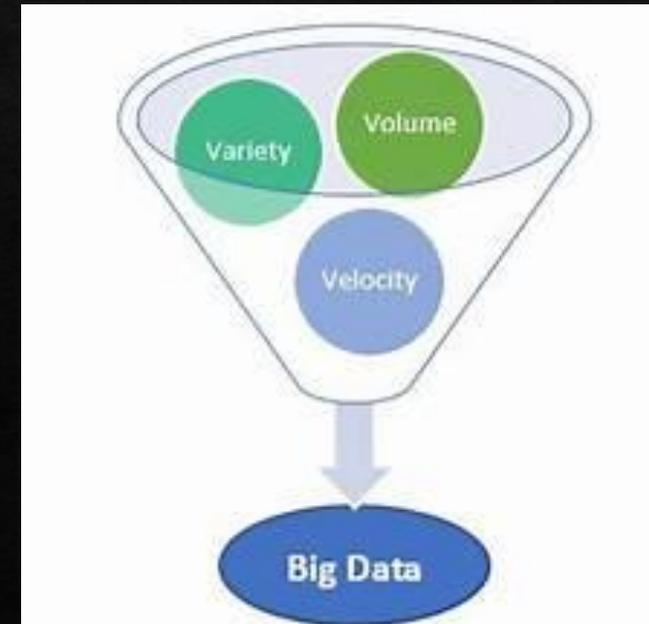


چالش‌های BigData

The 3 V's of Big Data



اصل 3V در
Big Data



90%
of today's data
has been created
in just the last
2 years



Every day
we create
2.5
quintillion
bytes of data

Every
60
seconds
there are

72 hours
of footage
uploaded to
YouTube

50,000
GB/second
is the estimated
rate of global
Internet
traffic
by 2018

216,000
Instagram posts

204,000,000
emails sent

(...enough to fill
10 million
Blu-ray
discs)

Volume
Scale of data

Velocity
Speed of data

Veracity
Certainty of data

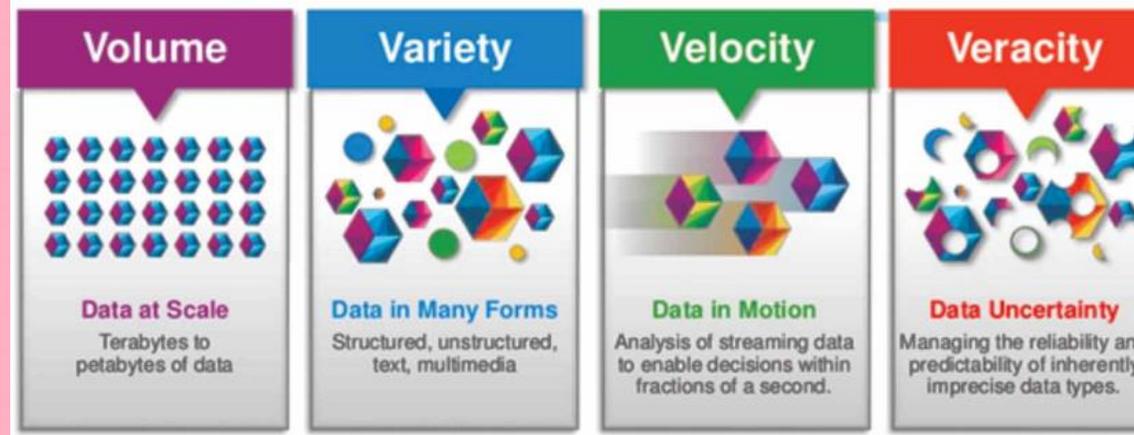
Variety
Diversity of data

1 in 3
business leaders
don't trust the
information they use
to make decisions

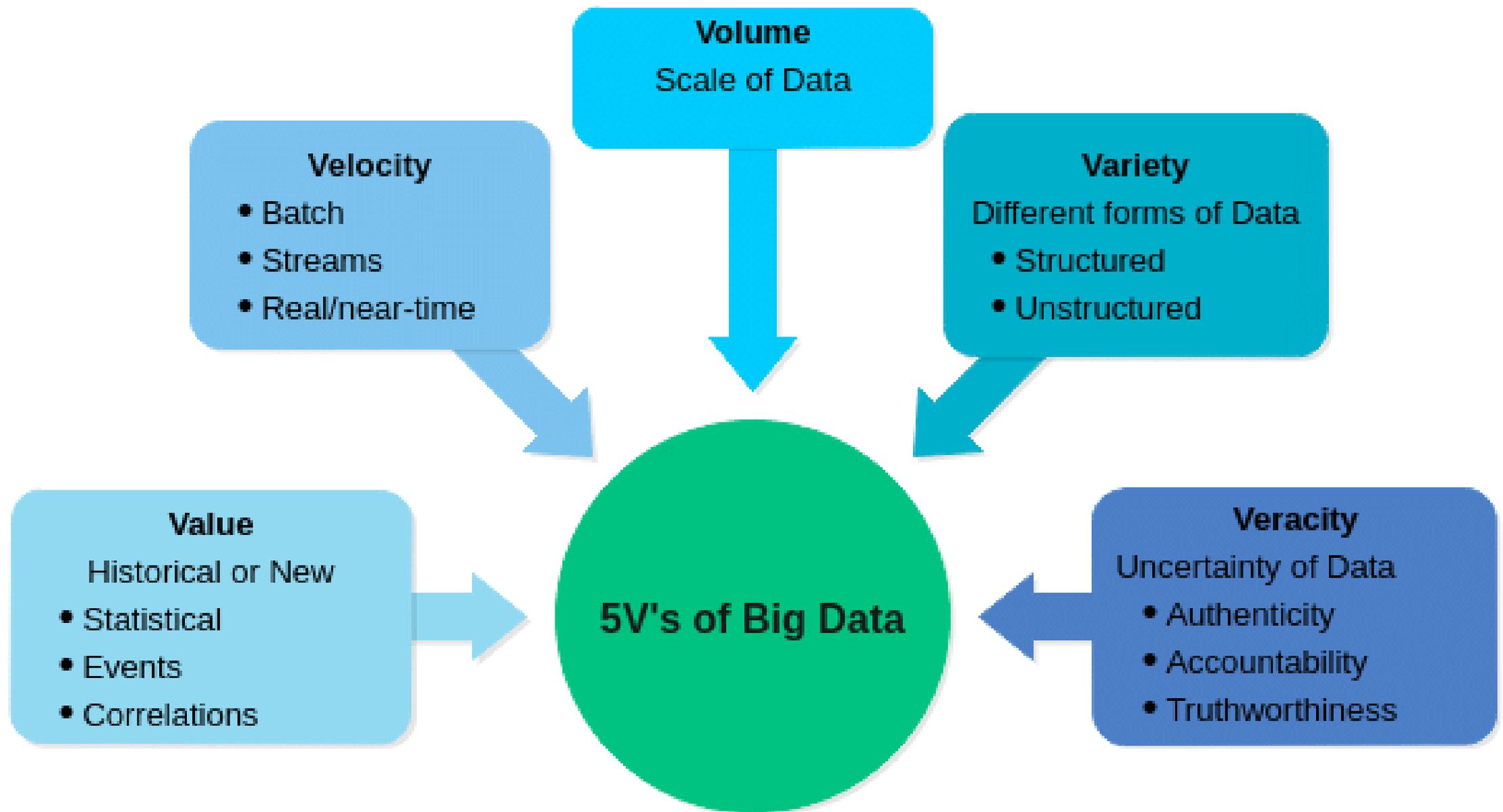
80% of data
growth is video,
images and
documents

\$3.1 trillion
is the estimated
amount of money that
poor data quality costs
the US economy per year

90%
of generated data
is "unstructured"
This includes tweets, photos,
customer purchase histories
and customer service calls



اصل 4V در Big Data





BIG DATA

with 8 V's

08

VIRALITY

AHA to-go? Does it convey a message that can be pasted into a presentation or Instagrammed?



01

VOLUME

Can you find the information you are looking for?



02

VALUE

Can you find it when you most need it?



03

VERACITY

Are you dealing with information or disinformation?



04

VISUALISATION

Can you make sense at a glance? Does it trigger a decision?



05

VARIETY

Is a picture worth a thousand words in 70 languages? Is your information balanced?



07

VISCOSITY

Does it stick with you? Does it call for action?



06

VELOCITY

Information gains momentum and crises & opportunities evolve in real time. How is outlook for today?



چالشهای کلان داده

The High Cost of Inaccurate Data Industry statistics reveals

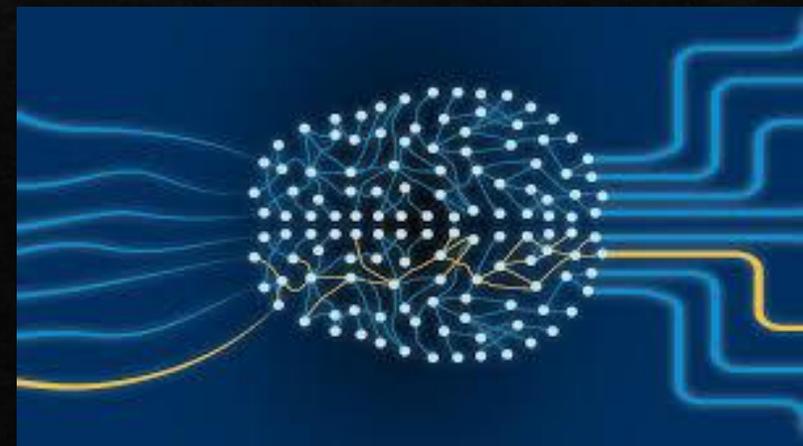
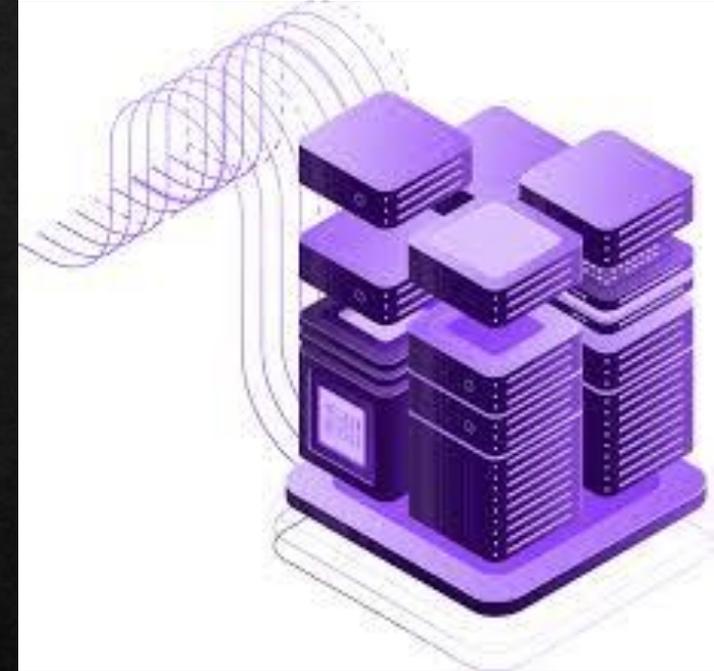
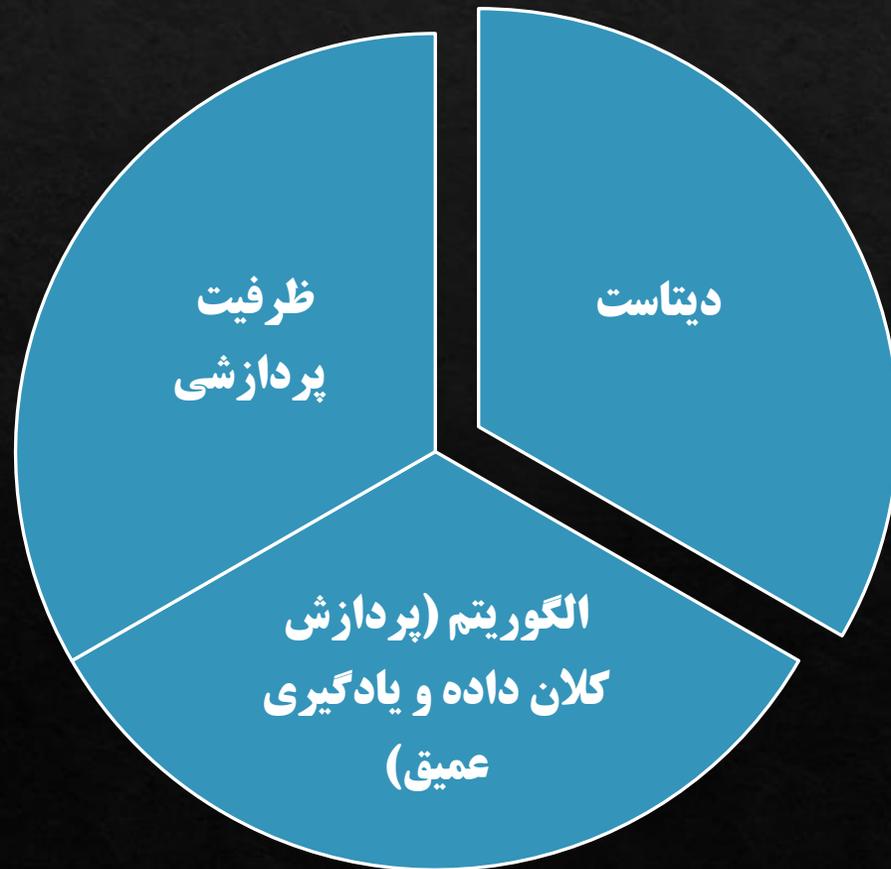


تحوالات جهانی بعد از حضور
کلان داده و هوش مصنوعی

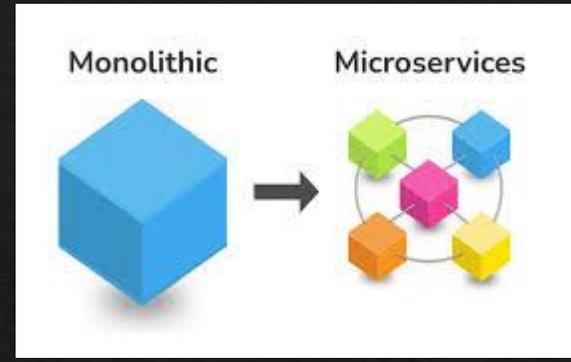
تحولات کلان داده در علوم مختلف در جهان

- ◇ حذف پزشکی سنتی و جایگزینی آن با پزشکی شخصی (Personalized medicine)
- ◇ حذف داروسازی سنتی و ایجاد داروی ویژه هر فرد با استفاده از پردازش داده ژنتیک
- ◇ از بین رفتن دموکراسی
- ◇ حذف رانندگی سنتی و ایجاد خودرو هوشمند و جاده هوشمند و نقش مهم رگولاتوری داده در آینده
- ◇ حذف وکالت در آینده
- ◇ نقش مهم مهندسان داده در آینده

سه راس مهم هوش مصنوعی



Microservice Architecture



تکنولوژی های کلیدی

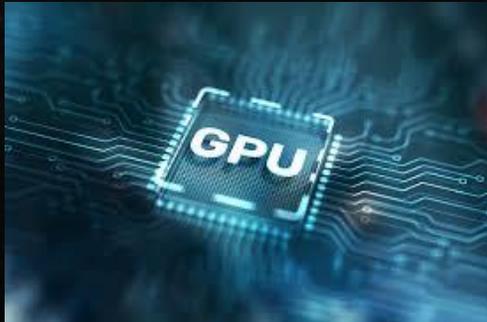
Kubernetes



Distributed Storage



Pipelining with Advanced GPU



استفاده از ۳۷ گیگابایت رم از ۴۰ گیگابایت رم A100 GPU فقط برای پردازش اطلاعات یک فرد

NVIDIA GH200

Built for the new era
of AI supercomputing

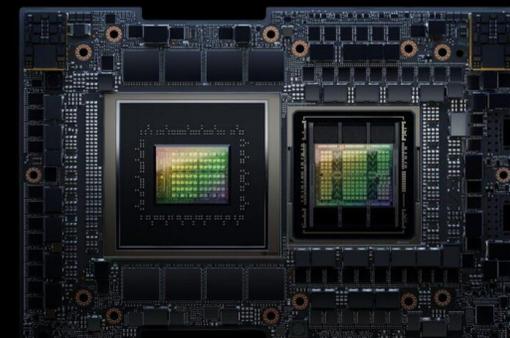
CPU to GPU Bandwidth
900 GB/s
NVLink-C2C

Memory Bandwidth
4.9 TB/s
HBM3e per GPU

Energy Efficiency
1.9X
Performance vs H100

QFT Quantum Simulation
90X
Performance vs dual x86 CPU

RAG LLM Inference
100X
Performance vs dual x86 CPU

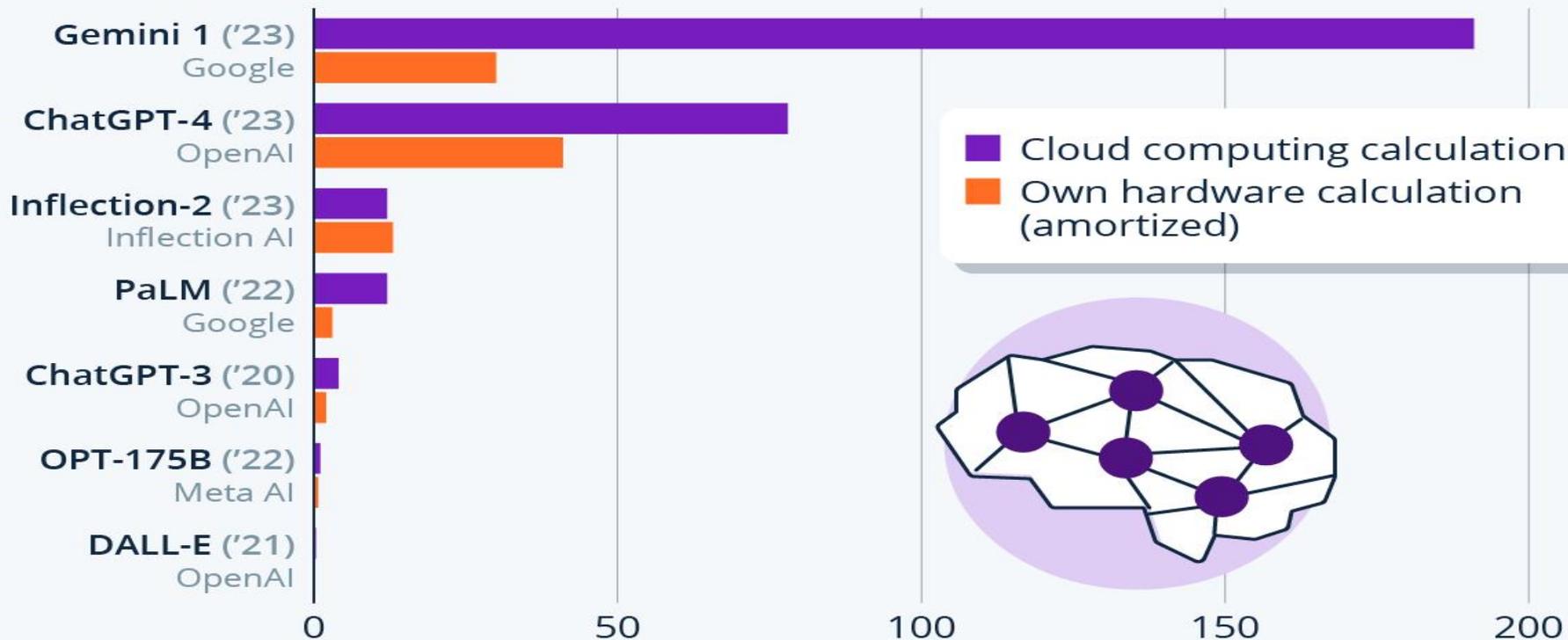


624 GB High-Speed Memory | 4.9 TB/s | 4 PF AI Perf | 72 Arm Cores

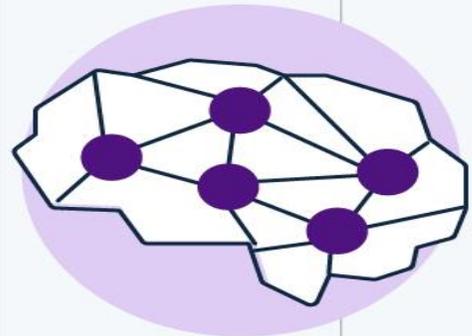


The Extreme Cost of Training AI Models

Estimated cost of training selected AI models (in million U.S. dollars), by different calculation models



Legend:
■ Cloud computing calculation
■ Own hardware calculation (amortized)



چالش‌های هزینه آموزش مدل

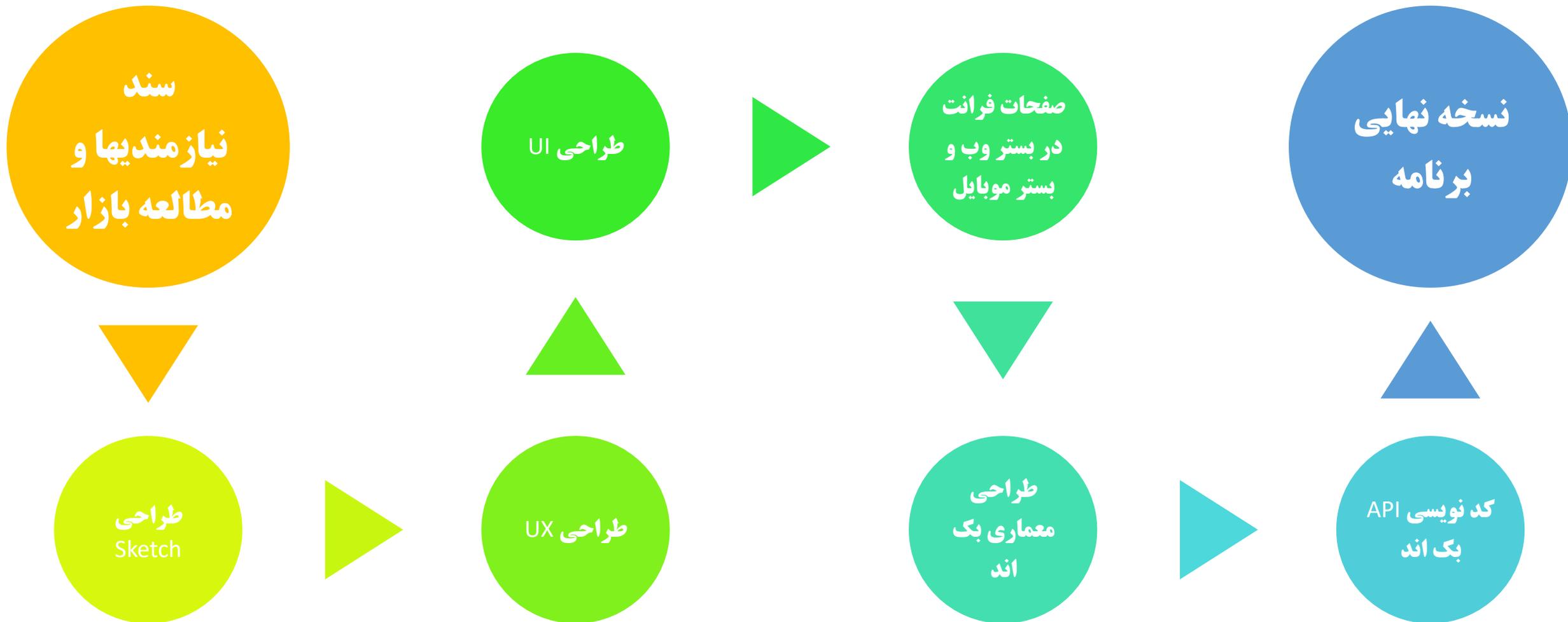
Rounded numbers. Excludes staff salaries that can make up 29-49% of final cost (including equity)

Source: Epoch AI



مراحل توسعه یک نرم افزار

مراحل توسعه نرم افزار



نمونه ای از صفحات App

در دستر Figma

MedicMetaVerse

File Assets

Pages

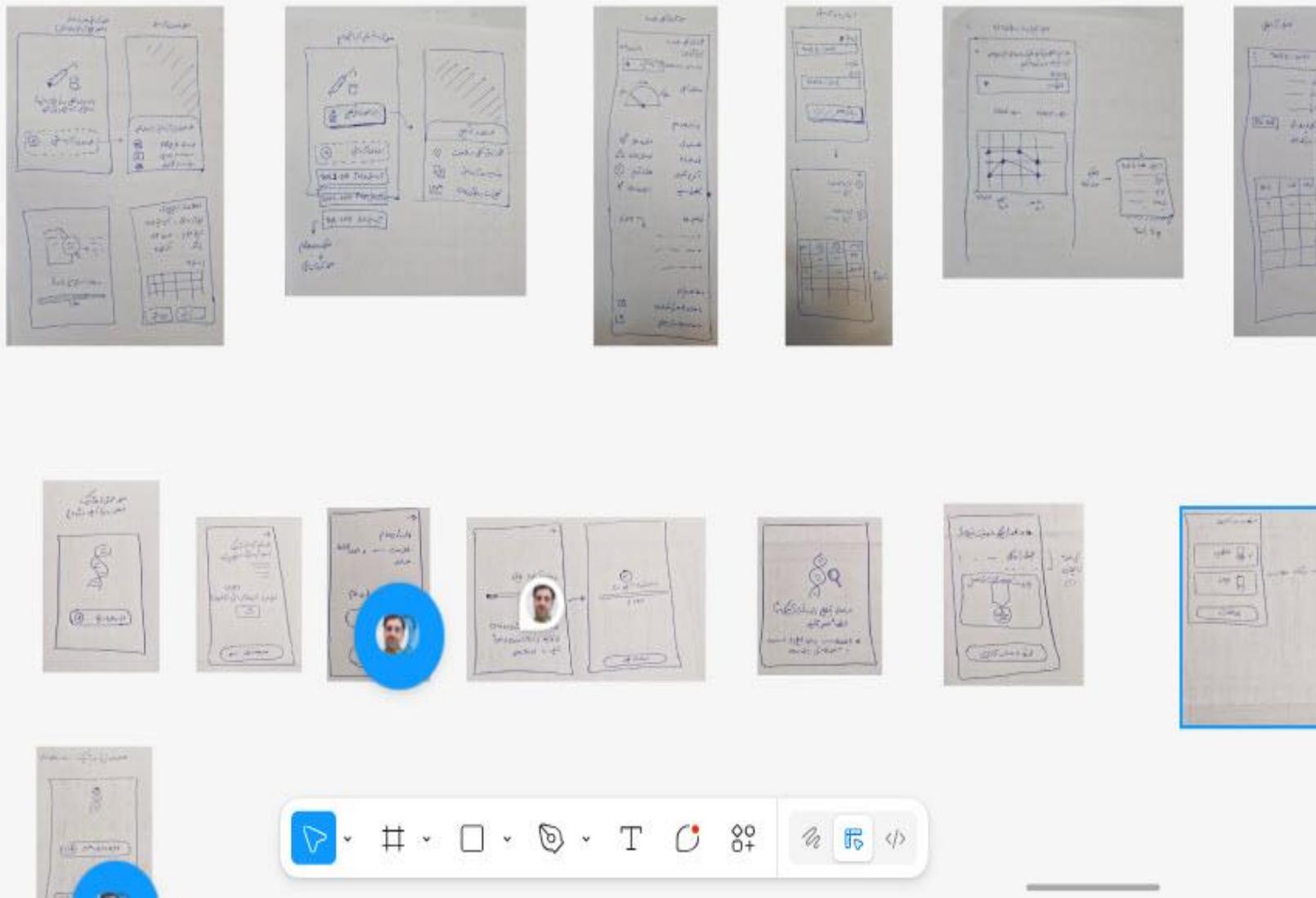
- Final
- Page 3
- Draft

Layers

- Frame 1437254888
- Frame 1437254887
- Anatomy
- image 48
- UI Kit
- Risks
- Assistant
- buy pack
- Incomplete Profile
- buy pack
- Incomplete Profile
- buy pack

Frame 1437254793

Sketch



Design Prototype 16%

Page

F6F6F6 100%

Variables

Styles

Text styles

- T 1 عنوان - 31/Auto
- T 2 عنوان - 25/Auto
- T 3 عنوان - 20/Auto
- T 1 بنه - 16/Auto
- T 2 بنه - 13/Auto
- T 3 بنه - 10/Auto

Export

Navigation bar with icons for selection, zoom, text, and other tools.

<https://www.instagram.com/aspiradesign/>

Miller's Law

The average person can only keep 7
(plus or minus 2) items in their working memory.



مقدمه

الگوهای برنامه‌نویسی

انواع دسته‌بندی زبان‌های
برنامه‌نویسی

انواع حوزه‌های کاری
برنامه‌نویسی

مقایسه‌ی زبان‌های
برنامه‌نویسی مختلف



مقدمه

الگوهای برنامه‌نویسی

انواع دسته‌بندی زبان‌های
برنامه‌نویسیانواع حوزه‌های کاری
برنامه‌نویسیمقایسه‌ی زبان‌های
برنامه‌نویسی مختلف

حالت‌های مختلف Input

برای طراحی ورودی‌ها (Input) چهار حالت مختلف که برای کاربر اتفاق می‌افتد را باید در نظر بگیریم.

شماره موبایل	حالت عادی
<input type="text" value="۹۱۲۱۲۳۴۵۶۷"/>	
شماره موبایل	حالت فعال (Active)
<input type="text"/>	
شماره موبایل	حالت پر (Filled)
<input type="text" value="۹۱۲۱۲۳۴۵۶۷"/>	
شماره موبایل	حالت خطا
<input type="text" value="۹۱۲۱۲۳۴۵"/>	

شماره موبایل وارد شده اشتباه است

انواع دکمه‌ها در طراحی

مدل‌های مختلف آیکن، دکمه به معنای یک عملیات (یا یک سری عملیات) است. با کلیک بر روی یک دکمه، یک عملیات اجرا می‌شود.

دکمه متنی	خرید از گیک‌شو
دکمه پیش‌فرض (Default)	<input type="button" value="خرید از گیک‌شو"/>
دکمه نقطه چین (Dashed)	<input type="button" value="خرید از گیک‌شو"/>
دکمه اصلی (Primary)	<input type="button" value="خرید از گیک‌شو"/>
دکمه غیر فعال (Disabled)	<input type="button" value="خرید از گیک‌شو"/>



مقدمه

الگوهای برنامه‌نویسی

انواع دسته‌بندی زبان‌های
برنامه‌نویسی

انواع حوزه‌های کاری
برنامه‌نویسی

مقایسه‌ی زبان‌های
برنامه‌نویسی مختلف

<https://www.instagram.com/wwuiuxdesign/>

Color Psychology Cheatsheet

Red

passion, love, anger,
danger, warning

Orange

warmth, energy,
happiness, enthusiasm

Green

nature, growth, health,
harmony, money

Yellow

sunshine, happiness, joy,
intellect, caution

Blue

calmness, trust, loyalty,
wisdom, mystery

Purple

royalty, luxury, creativity,
mystery

Pink

innocence, femininity,
compassion, love

Black

power, sophistication,
mystery, evil

White

purity, innocence,
cleanliness, peace

Brown

earthiness, reliability,
dependability, friendliness

Silver

modernity, sophistication,
industrial, cold

Gold

luxury, wealth, success,
prestige, glamor



مقدمه

الگوهای برنامه‌نویسی

انواع دسته‌بندی زبان‌های برنامه‌نویسی

انواع حوزه‌های کاری برنامه‌نویسی

مقایسه‌ی زبان‌های برنامه‌نویسی مختلف

01. Colors

| COLORS

PRIMARY COLORS



ACCENT COLORS

NEUTRAL COLORS



02. Typography

Raleway

Heading 1

Raleway SB
34px L41

Heading 2

Raleway SB
22px L28

Heading 3

Raleway SB
17px L22

Heading 4

Raleway M
15px L20

03. Components

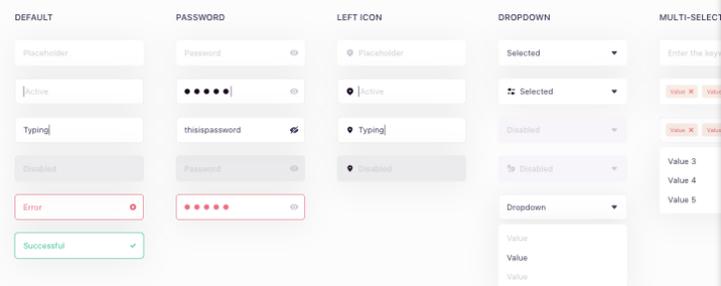
| COMPONENTS

BUTTONS

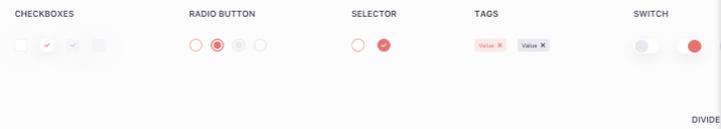
LARGE BUTTONS



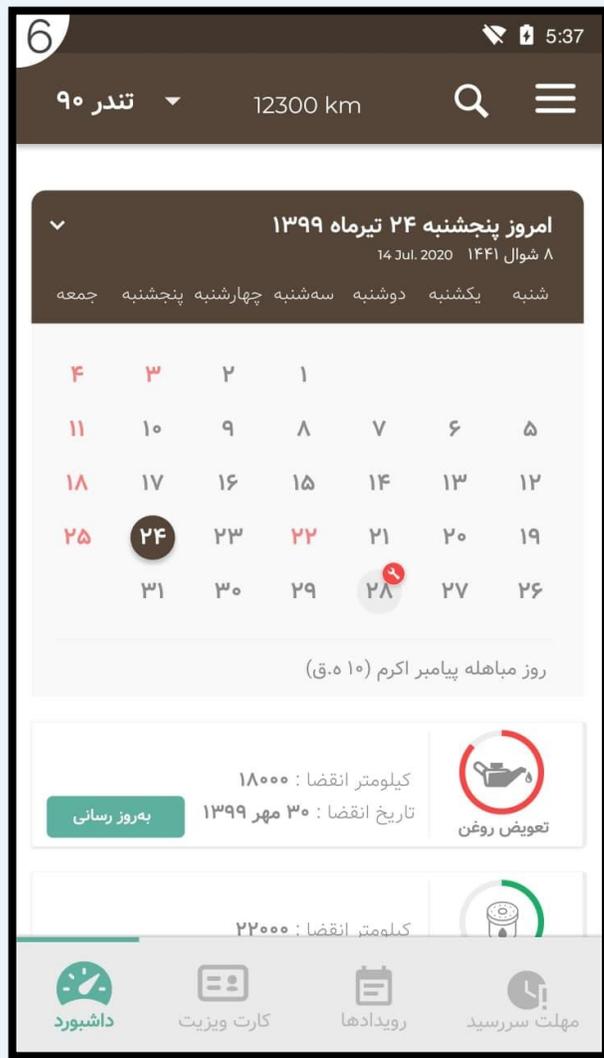
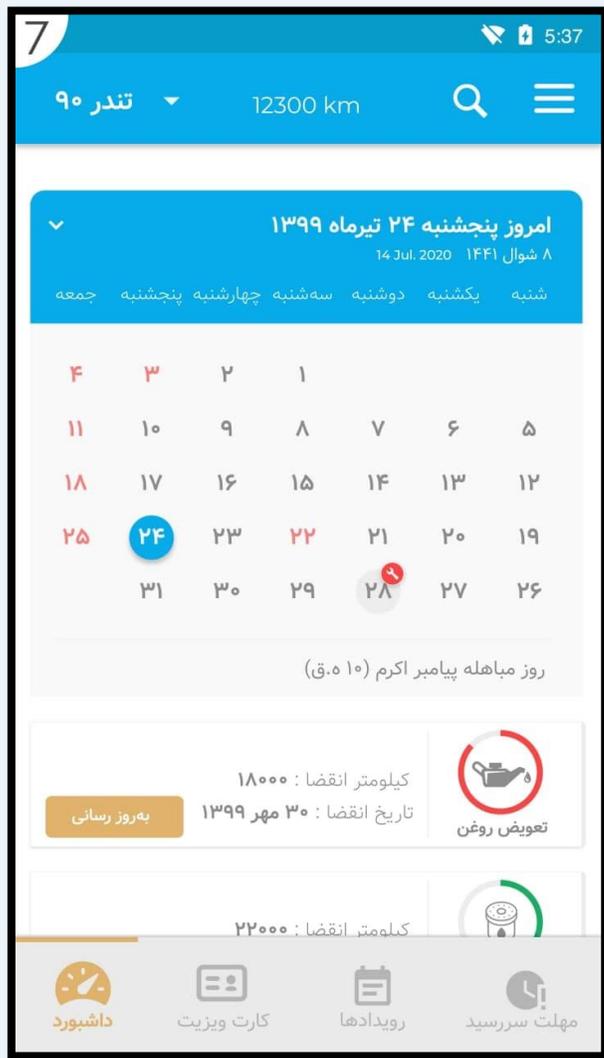
INPUT FORMS



OTHER COMPONENTS



Primary Color



مقدمه

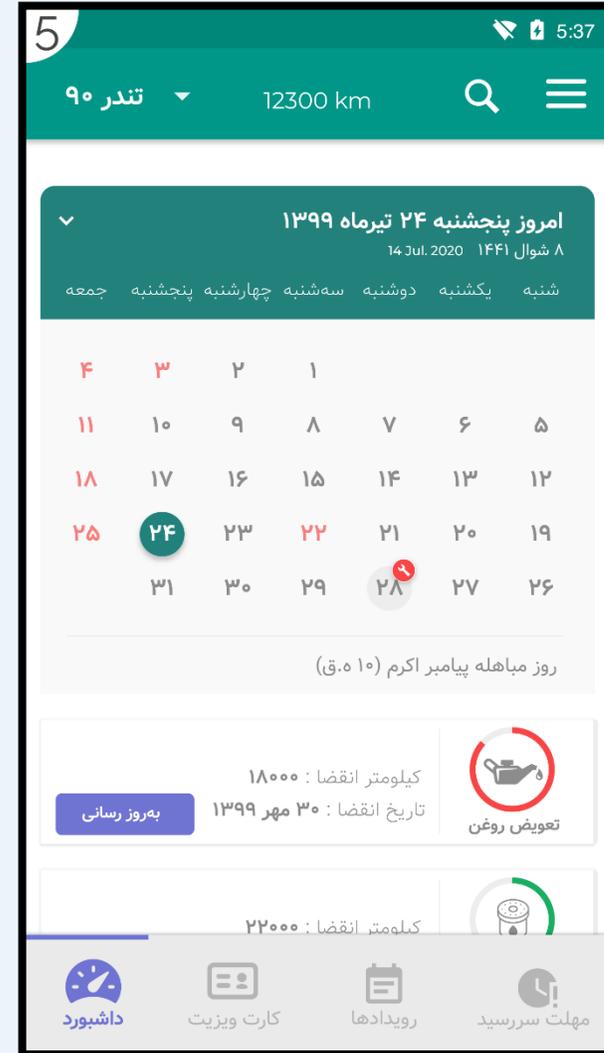
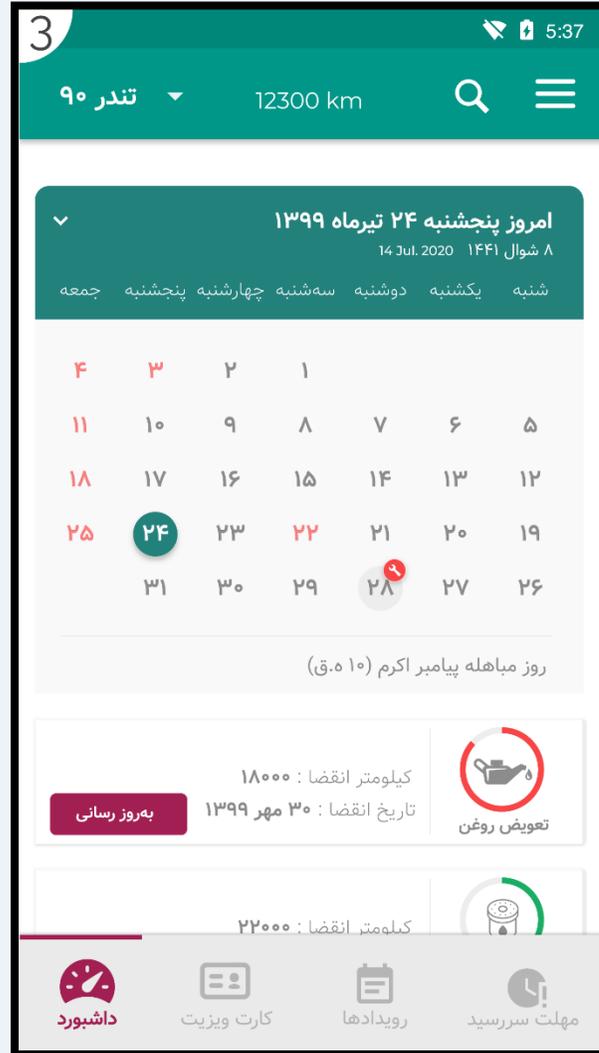
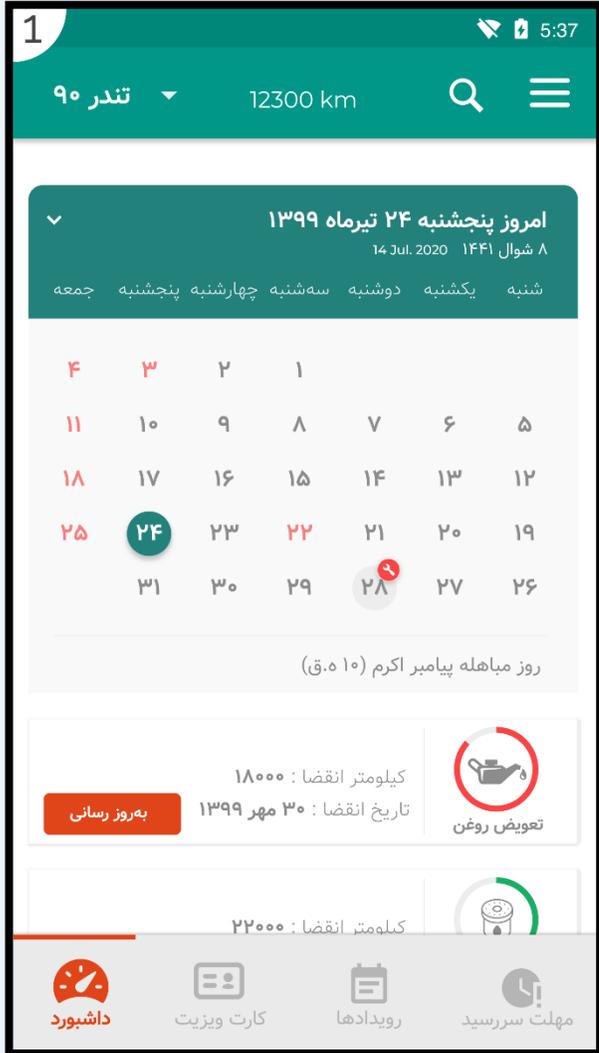
الگوهای برنامه‌نویسی

انواع دسته‌بندی زبان‌های برنامه‌نویسی

انواع حوزه‌های کاری برنامه‌نویسی

مقایسه‌ی زبان‌های برنامه‌نویسی مختلف

Secondary Color



مقدمه

الگوهای برنامه‌نویسی

انواع دسته‌بندی زبان‌های برنامه‌نویسی

انواع حوزه‌های کاری برنامه‌نویسی

مقایسه‌ی زبان‌های برنامه‌نویسی مختلف

Colors



Typography

	عنوان	العنوان	نص
عنوان 1	16pt	16pt	14pt
عنوان 2	14pt	14pt	12pt
عنوان 3	12pt	12pt	10pt
...	10pt	10pt	8pt
...	8pt	8pt	6pt
...	6pt	6pt	4pt

BG



iPhone X...



Bu...



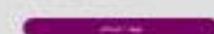
B...



C...



Action Bar



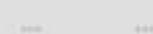
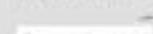
InputFie...



Main-He...



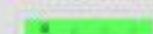
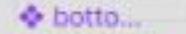
InputFie...



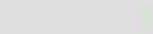
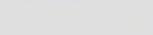
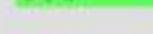
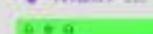
navitem



botto...



Main-...



UI Kit



توصیه‌های دستیار پزشک

بر اساس وضعیت کلی سلامت بهتر است موارد زیر را رعایت کنید و در صورت لزوم به پزشک مراجعه کنید

- مصرف چربی‌های ترانس را به حداقل برسانید
- هر شش ماه یکبار چکاپ آزمایش‌های کبدی را تکرار کنید
- مصرف سبزیجات روزانه را در رژیم خود بگنجانید
- سریع‌تر آزمایش تری‌گلیسرید را انجام دهید

مشاهده همه توصیه‌ها ←

بیشتر بخوانید

مقالات زیر کمک می‌کنند اطلاعات بیشتری در مورد سلامت خود داشته باشید

[از آزمایش خون چه می‌دانیم؟](#)

[فرق کلسترول خوب و کلسترول بد چیست؟](#)

[آیا برای آزمایش ادرار باید ناشتا باشیم؟](#)

وضعیت ژنتیکی

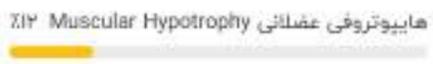
ژن	شرایط	تاثیر
BRCA1	Breast Cancer	High
MGF-3	Muscle Degradation	Low
BRCA1	Breast Cancer	Medium
IL-13RA	Asthma	Low
BRCA1	Breast Cancer	High
BRCA1	Breast Cancer	High

مشاهده همه چشم‌ها ←



ریسک بیماری‌ها

نقشه ژنتیک شما نشان‌دهنده ریسک ابتلا به برخی بیماری‌ها است. لیست زیر بر اساس میزان ریسک بیماری‌ها مرتب شده است



مشاهده همه ریسک‌ها ←



وضعیت کلی سلامت



← آزمایش خون و ادرار
بسته فعال ندارید

← آزمایش ژنتیک
نمونه اعتبار بسته تا ۲۴ اکتبر ۱۴۰۴

← مشاهده همه امکانات

وبلاگ

نقشه‌برداری ژنی در قرن ۲۱ چه پیشرفت‌هایی داشته؟
۴ مرداد ۱۴۰۲ - ژنتیک
هوش مصنوعی به کمک هئندسین ژنتیک آمده و نقشه‌برداری ژنتیکی بسیار دقیق...



تاثیر ژنتیک بر ورزش
۲۰ شهریور ۱۴۰۲ - عمومی - ژنتیک
ورزشکاران نظیرام نورذرات فرزان گاهی در حصار ژنهای خود اسیر هستند و بیش...



روش جدید نمونه‌برداری خونی
۱ فروردین ۱۴۰۲ - علوم آزمایشگاهی
بدون درد هوش مصنوعی به کمک هئندسین ژنتیک آمده و نقشه‌برداری ژنتی...



← مشاهده همه مطالب





بررسی و تحلیل آزمایش

آزمایش خون ۱۰۱



آزمایشگاه سلامت فردا
دکتر علی محمدی
۲۳ آذر ۱۴۰۳

افزودن آزمایش

آزمایش با موفقیت ثبت شد

نوع آزمایش: خون

تاریخ انجام آزمایش: ۲۳ آذر ۱۴۰۳

نام پزشک: دکتر علی محمدی

آزمایشگاه: سلامت فردا

آدرس آزمایشگاه: قلعه کوچک اسفندیاری

شماره آزمایش: ۳۸۹۱۰۳

فایل آپلود شده: 20250312221234.PDF

پارامترهای استخراج شده

با کلیک روی هر سطر راهنمای مربوط به پارامتر را ببینید

پارامتر	واحد	مقدار	حد
WBC	#/ml	12000	00
RBC	#/ml	100000	00
Plaquette	#/ml	100000	00
pH	-	6.9	۱.2
Triglyceride	mg/l	100000	00
RBC	#/ml	100000	50
Cholesterol	mg/l	12000	00

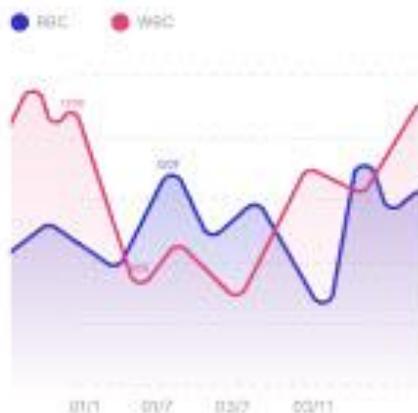
تمام صفحه

ذخیره

اشتراک‌گذاری

چند پارامتر را انتخاب کنید

پارامترهای آزمایش را انتخاب کنید و تغییرات زمانی را در چارت زیر مشاهده کنید.



نمایش اعداد روی چارت

RBC

WBC

افزودن پارامتر

ذخیره چارت

اشتراک‌گذاری

نوع فایل: Whole Genome

تاریخ انجام آزمایش: ۲۳ آذر ۱۴۰۳

آزمایشگاه: سلامت فردا

شماره آزمایش: IPA۰۱۰۱۰۳



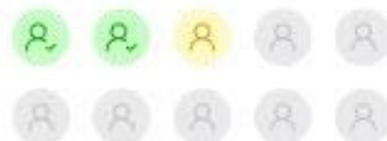
با خرید بسته جمعیت اجزای خود تباریابی ژنتیکی دقیقتر خود را در کشورهای مختلف روی نقشه ببینید

[بچشم مشاهده نتایج کامل](#)


تعداد	درصد
آسیایی	۷۴۵
قفقازی	۷۱۲
ایرانی	۷۱۱
ترکیه‌ای	<۷۱
افغان	<۷۱
اروپایی	۷۳۵
ایتالیایی	۷۲۵
اسپانیایی	۷۷
فرانسوی	۷۳
آفریقایی	۷۲۰
مصری	۷۱۸
الجزایری	۷۲

با دعوت دوستان خود بسته رایگان هدیه بگیرید

با دعوت ۵ نفر دیگر به برنامه می‌توانید یک بسته رایگان آزمایش خون یا ژنتیک به انتخاب خود هدیه بگیرید.



۷ نفر دیگر تا دریافت هدیه

- ✓ علیرضا محمدی ثبت‌نام کرد
- ✓ بهرام سعیدی بسته برنامه را خریداری کرد
- ✓ محمد علوی بسته برنامه را خریداری کرد

لینک زیر را برای دوستان خود بفرستید

لینک زیر را کپی کنید و برای دیگران بفرستید یا این‌که از آن‌ها بخواهید کد QR زیر را اسکن کنند

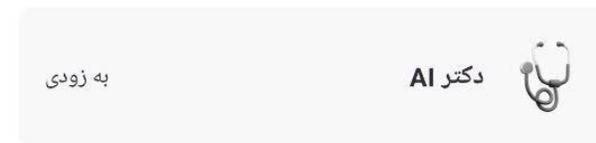
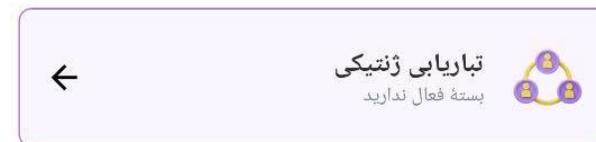
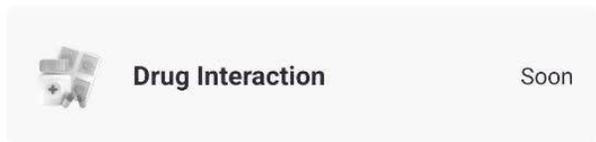
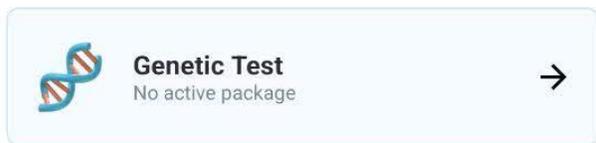
لینک اختصاصی شما

Medic.com/referral=12243908

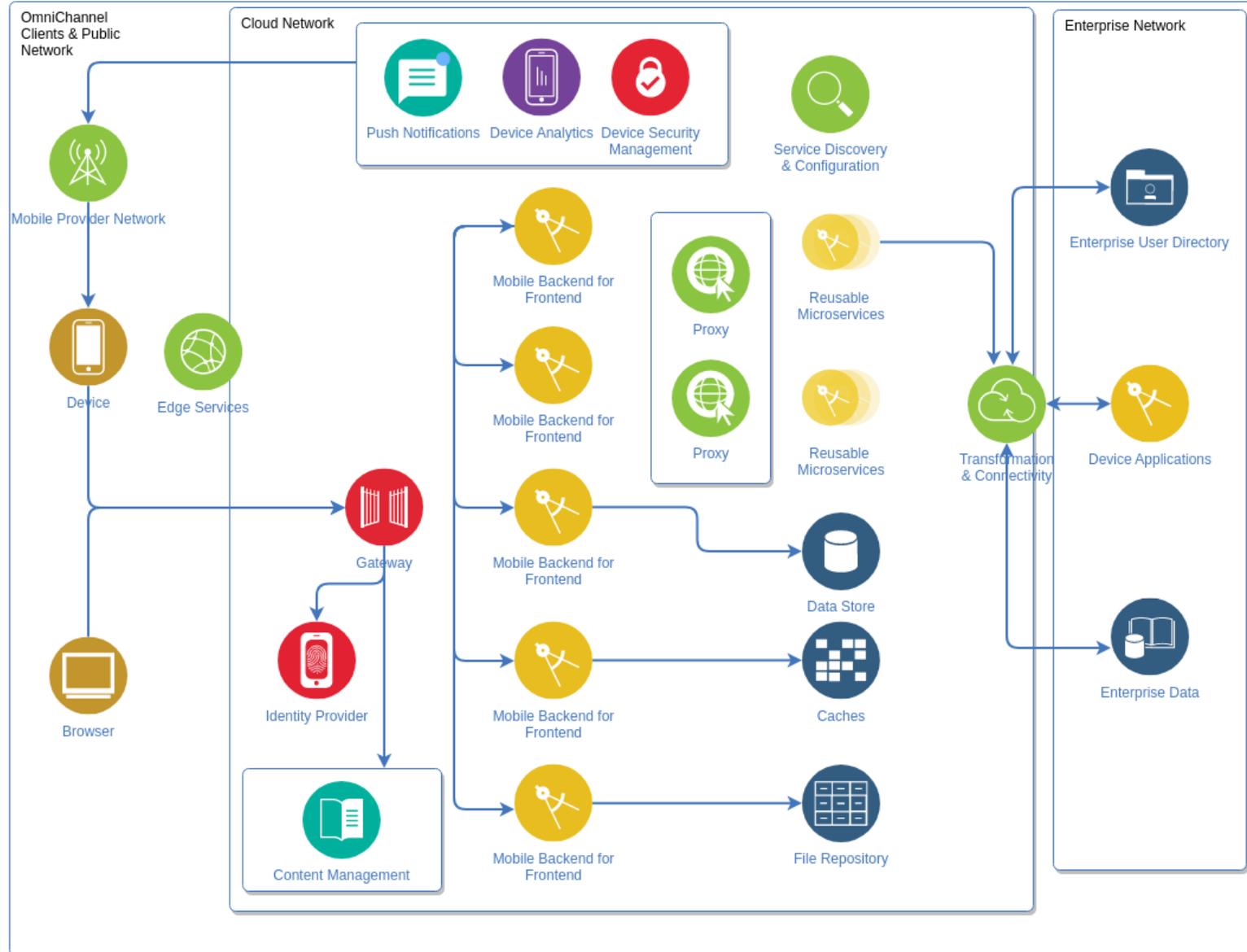
کپی لینک



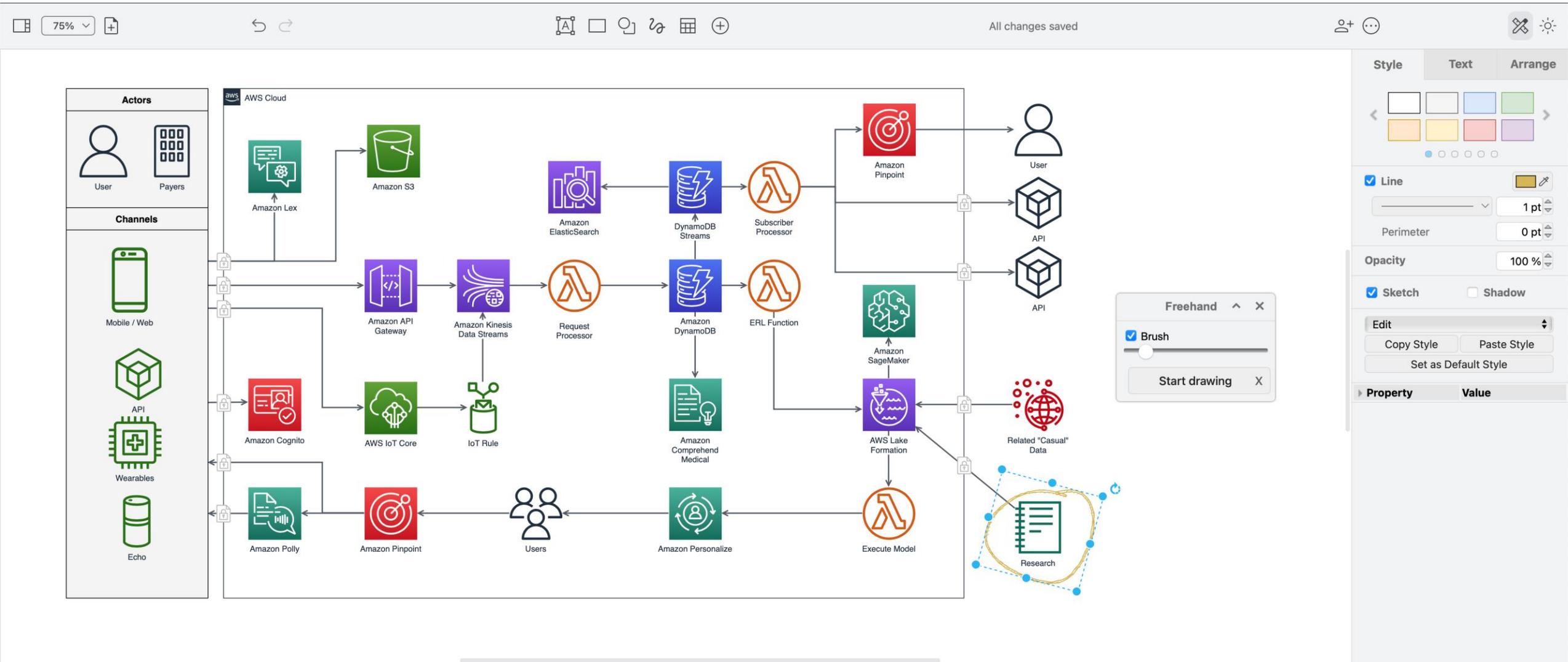
زبان های فارسی، عربی، انگلیسی



معماری میکروسرویس برنامه در بخش بک اند

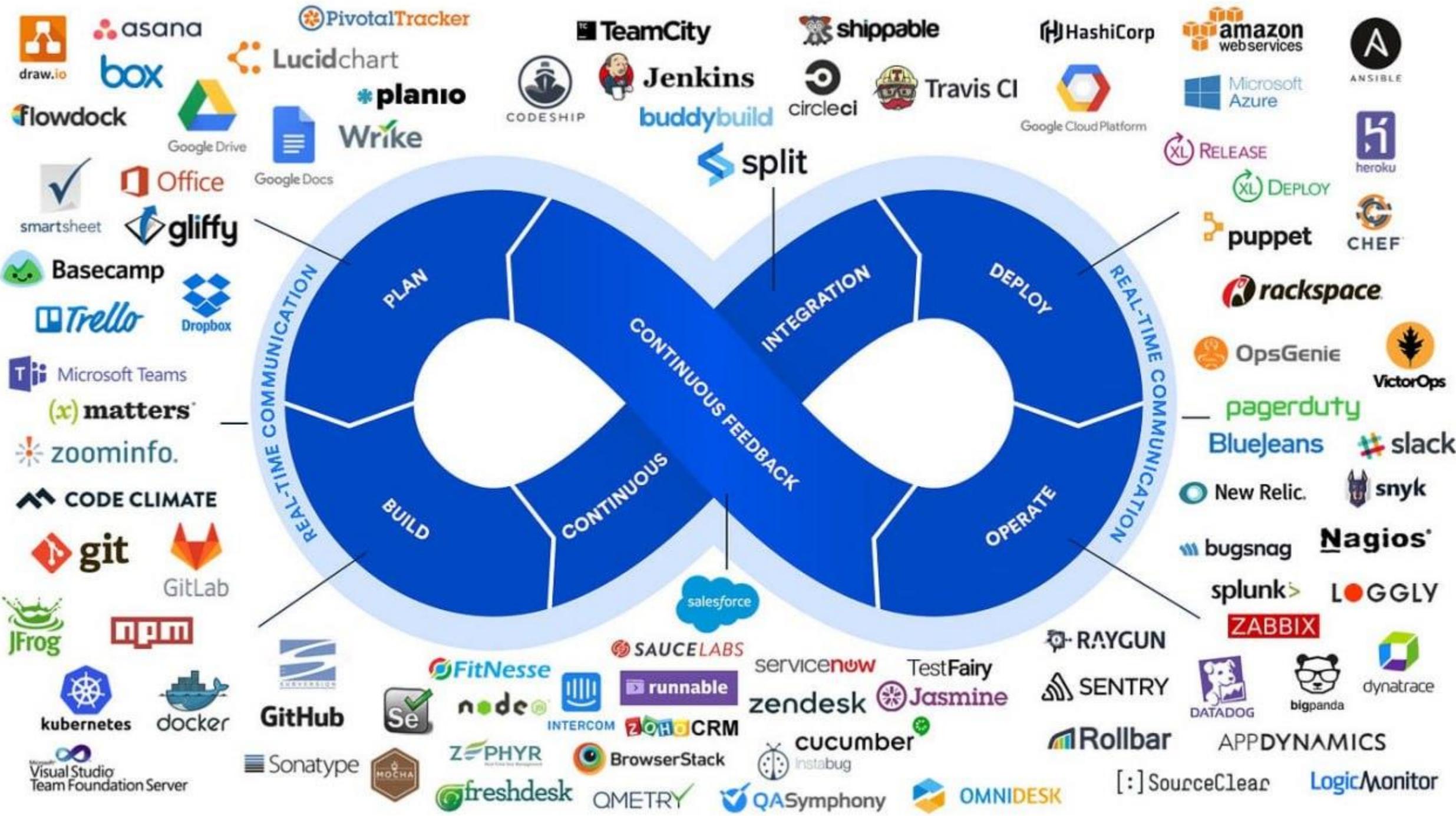


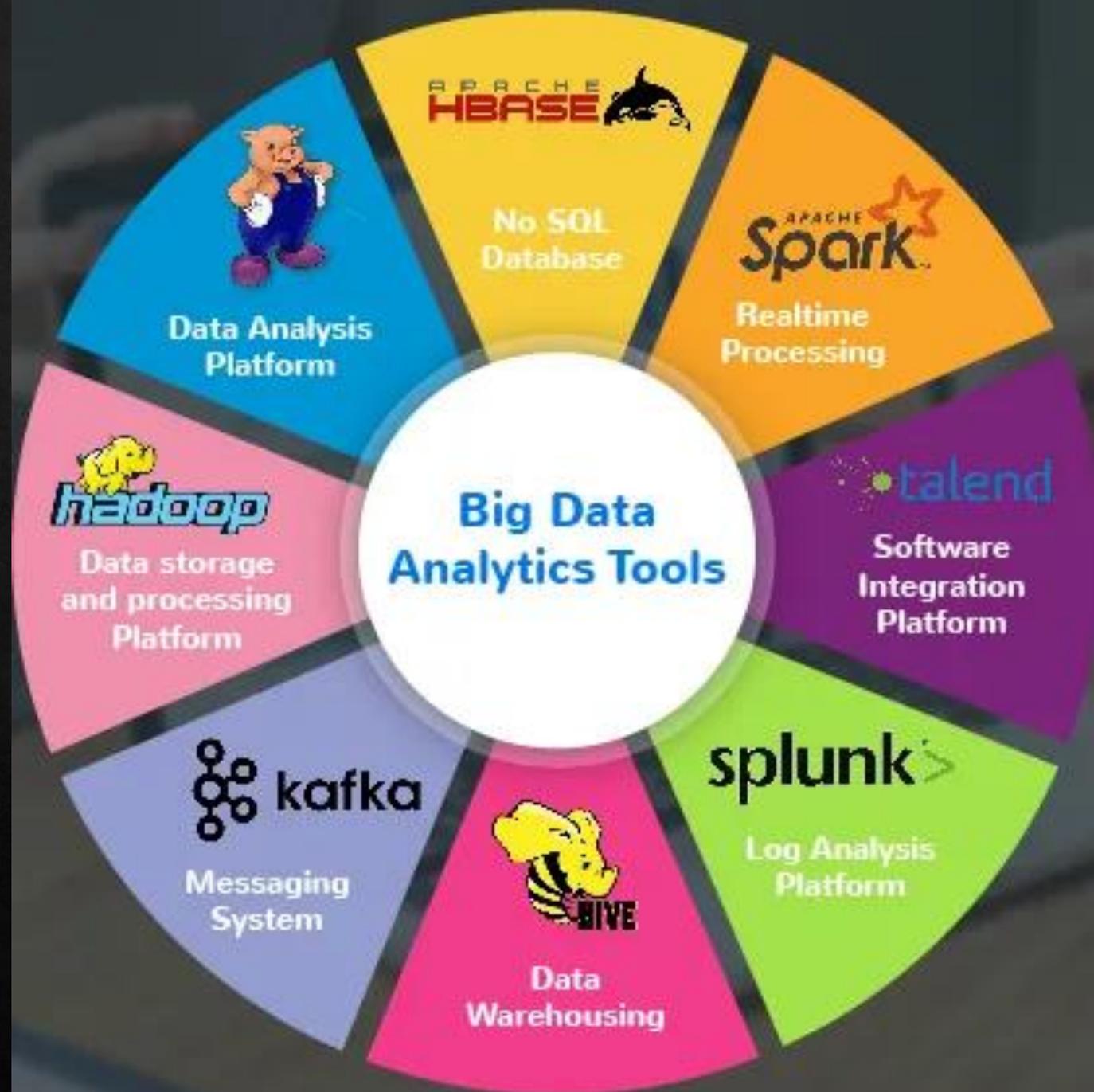
معماری Devops



نقش مهم ابزارها در حوزه برنامه نویسی و

Devops





Big Data Analytics Tools.



Open-source cross-platform source that utilizes a document oriented program.



An open-source platform that helps in the distribution and storage of large data.



Offers a larger insight into the hypothesis generated.



Helps in analyzing and manipulating the information through the use of visual programming.



A free open-source network analysis and visualization software tool.



Microsoft HDInsight is a big data solution powered by Apache Hadoop.



NoSQL databases are used to store unstructured data which have no particular scheme.



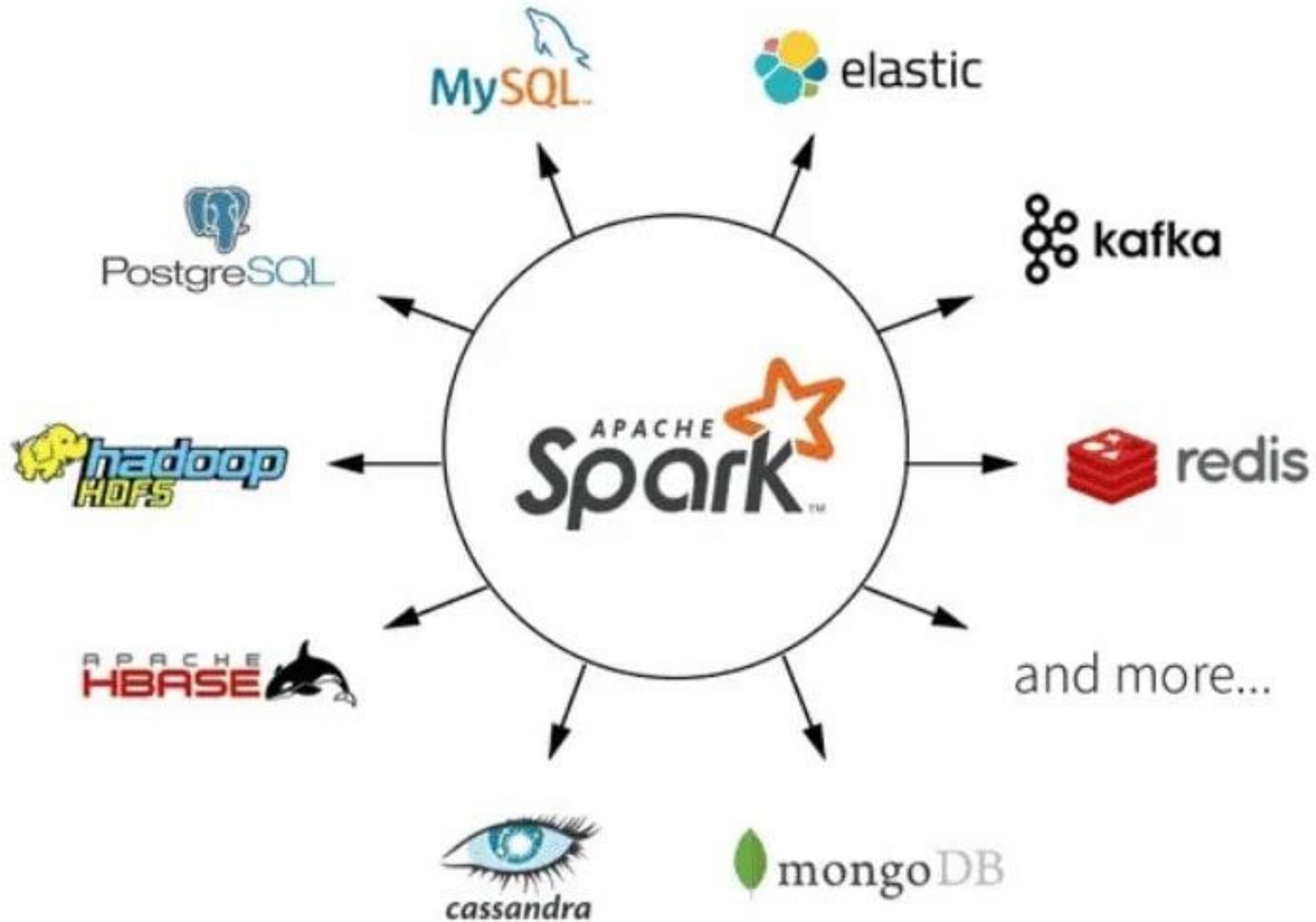
A distributed data management for Hadoop used for Data mining purpose.



A tool that is used to connect Hadoop with various relational databases that are used to transfer data.

Five Great Big Data Tools





مہارت ہا می مورد نیاز

Machine Learning

- Classification
- Regression
- Reinforcement Learning
- Deep Learning
- Clustering
- Dimensionally reduction

Programming Language

- Python
- R
- Java

Data Visualization

- Tableau
- Power BI
- Matplotlib
- GG Plot
- Seaborn

Data Analysis

- Feature Engineering
- Data Wrangling
- EDA

Data Science



IDE

- Pycharm
- Jupyter
- Colaboratory
- Spyder
- R-Studio

Math

- Statistics
- Linear Algebra
- Differential Calculas

Deploy

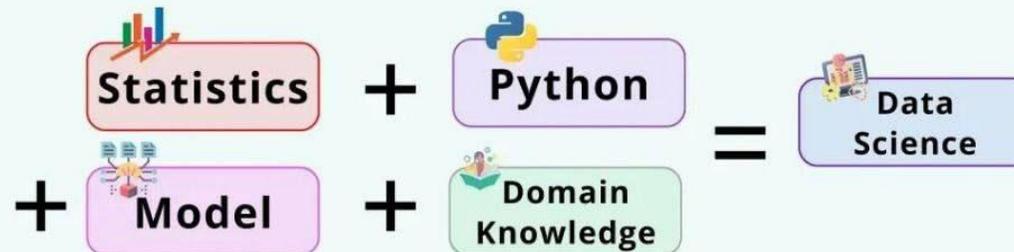
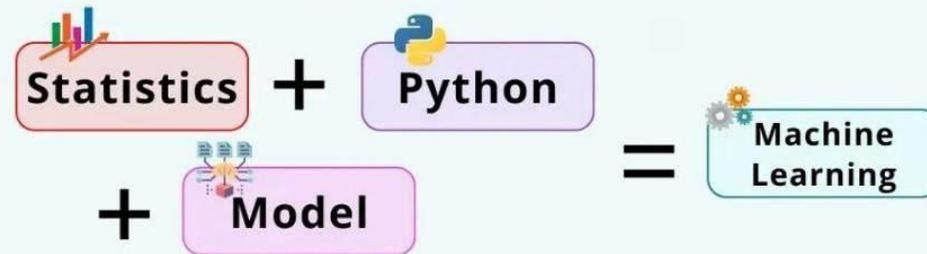
- AWS
- AZURE

Web Scraping

- Beautiful Soup
- Scrapy
- URLLIB

DATA SCIENCE

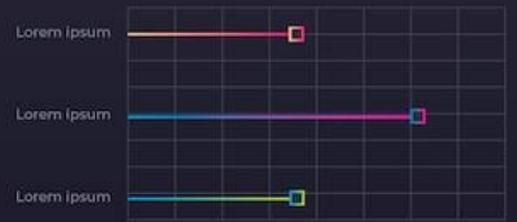
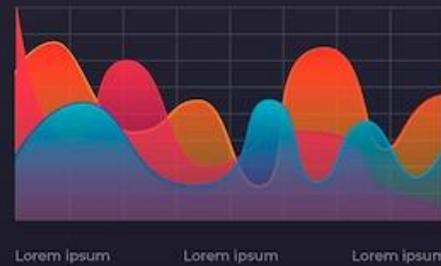
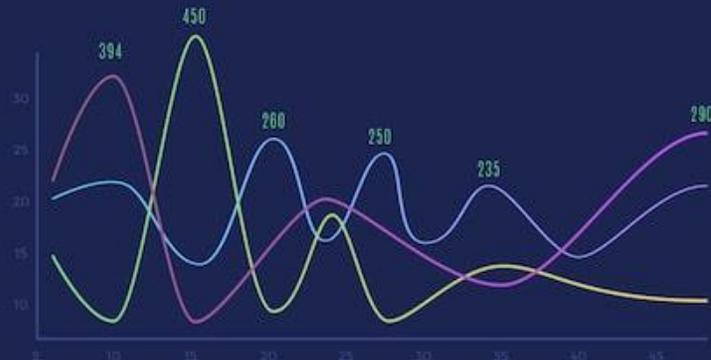
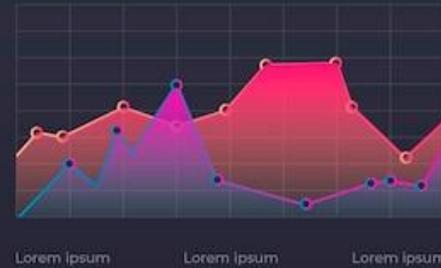
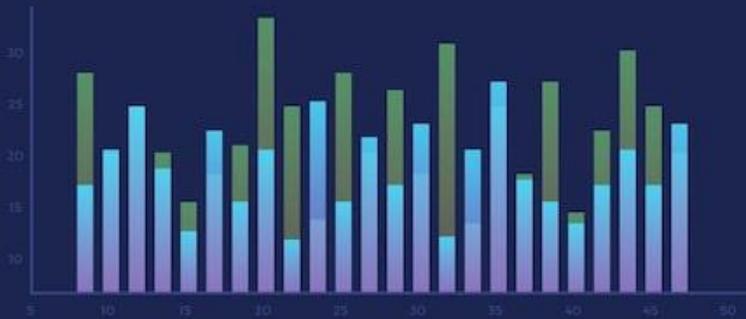
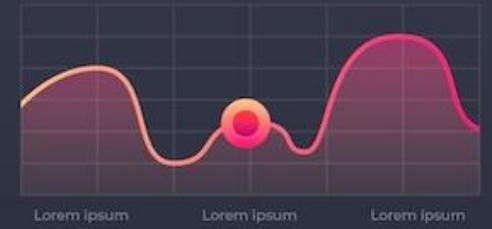
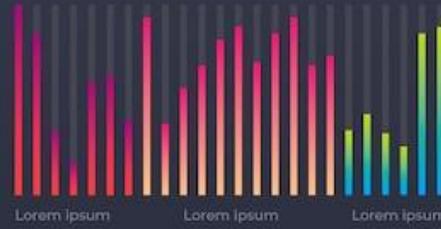
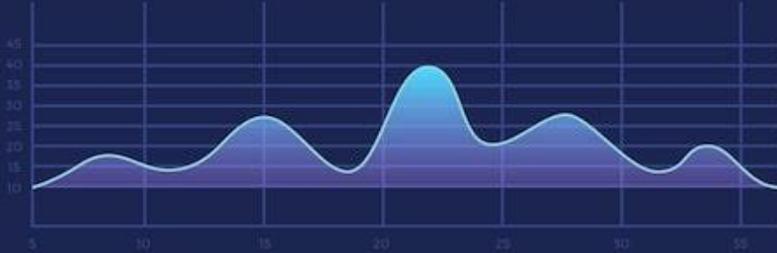
By: @pythoncodess



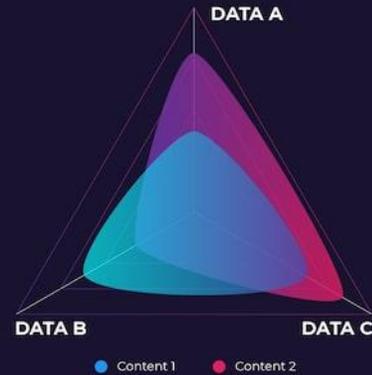
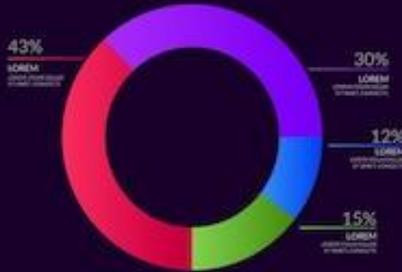
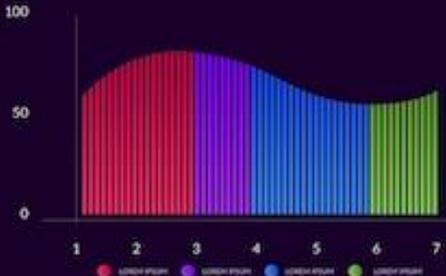
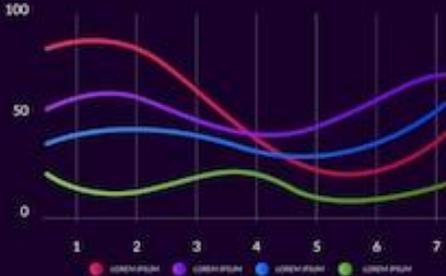
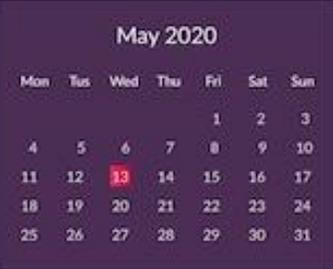
نمایش داده ها

Visualization

Visualization of data



Visualization of data



کتاب مرتبط با درس

بیشترین کتب مورد استفاده دانشگاههای برتر جهان در درس کلان داده

- ◇ *Database Systems: The Complete Book* — Garcia-Molina, Ullman, Widom web.cs.wpi.edu
- ◇ *Modern Database Management* — Hoffer, Ramesh, Topi web.cs.wpi.edu
- ◇ *Mining of Massive Datasets* — Leskovec, Rajaraman, Ullman library.buid.ac.ae
- ◇ *Hadoop: The Definitive Guide* — Tom White cs.sjtu.edu.cn
- ◇ *Learning Spark* — Holden Karau et al.

کتاب ۱۰ دانشگاه رتبه برتر دنیا در کلان داده

مثال‌هایی از دانشگاه‌ها و منابعی که استفاده می‌کنند 🔍

کتاب‌ها و منابع مورد استفاده

دانشگاه / دوره

از *Mining of Massive Datasets* — Jure Leskovec, Anand Rajaraman, Jeffrey Ullman استفاده می‌شود. mmds.org

Stanford University
دوره *Mining Massive Datasets (CS246)*

کتاب پیشنهادی شامل:

web.cs.wpi.edu *Database Systems: The Complete Book* — Garcia-Molina, Ullman & Widom •

web.cs.wpi.edu *Modern Database Management* — Hoffer, Ramesh, Topi •

web.cs.wpi.edu *Principles of Distributed Database Systems* — Tamer Ozsu & Patrick Valduriez •

Worcester Polytechnic Institute
(WPI)

دوره Big Data Management

eecs.yorku.ca کتاب اصلی: *Mining of Massive Datasets, 2nd Edition*

York University
دوره Big Data Systems

کتاب *Mathematics of Big Data: Spreadsheets, Databases, Matrices, and Graphs* — Jeremy Kepner & Hayden Jananathan در برخی دوره‌ها استفاده می‌شود. [2+ ...ate.of.Technology](http://2+...ate.of.Technology)

MIT

Examples of Universities and the References They Use in Big Data Courses

University / Course	Books and References Used
<p>Stanford University <i>Mining Massive Datasets (CS246)</i></p>	<p>Uses <i>Mining of Massive Datasets</i> — Jure Leskovec, Anand Rajaraman, Jeffrey Ullman (mmds.org ↗)</p>
<p>Worcester Polytechnic Institute (WPI) <i>Big Data Management</i></p>	<p>Suggested books include:</p> <ul style="list-style-type: none">• <i>Database Systems: The Complete Book</i> — Garcia-Molina, Ullman & Widom (web.cs.wpi.edu ↗)• <i>Modern Database Management</i> — Hoffer, Ramesh, Topi (web.cs.wpi.edu ↗)• <i>Principles of Distributed Database Systems</i> — Tamer Ozsu & Patrick Valduriez (web.cs.wpi.edu ↗)
<p>York University <i>Big Data Systems</i></p>	<p>Main textbook: <i>Mining of Massive Datasets, 2nd Edition</i> (eecs.yorku.ca ↗)</p>
<p>MIT</p>	<p>In some courses, they use <i>Mathematics of Big Data: Spreadsheets, Databases, Matrices, and Graphs</i> — Jeremy Kepner & Hayden Jananathan (mit.edu ↗)</p>

کتاب مفید درس Big data برای شروع و دید کلی

◇ برای شروع و دید کلی

◇ **Big Data: A Revolution That Will Transform How We Live, Work, and Think –**

◇ نوشته Kenneth Cukier و Viktor Mayer-Schönberger

کتابی غیر فنی که بیشتر درباره اهمیت داده‌های عظیم و تأثیرشان بر اقتصاد و جامعه است.

کتاب مفید درس Big data برای یادگیری مفاهیم فنی و ابزارها

◆ Hadoop: The Definitive Guide – ◆

◆ نوشته Tom White

مرجع کلاسیک برای یادگیری Hadoop زیرساخت پایه‌ای Big Data.

◆ Big Data: Principles and Best Practices of Scalable Real-Time Data Systems – ◆

◆ نوشته James Warren و Nathan Marz

کتابی عالی برای درک اصول طراحی سیستم‌های داده عظیم (از دید معماری).

◆ Streaming Systems – ◆

◆ نوشته Tyler Akidau و همکاران (مهندسان گوگل)

مرجع بسیار خوب برای پردازش داده‌های جریانی stream processing

کتاب مفید درسی Big data برای سطح پیشرفته و کاربردی

◆ **Designing Data-Intensive Applications** – نوشته Martin Kleppmann

شاید محبوب‌ترین کتاب حال حاضر در جامعه داده؛ مباحثی مثل معماری داده، سیستم‌های توزیع‌شده، پایگاه داده‌های NoSQL، پردازش real-time و batch را پوشش می‌دهد.

◆ **Mining of Massive Datasets** – نوشته Jure Leskovec, Anand Rajaraman, Jeffrey Ullman

کتابی دانشگاهی با الگوریتم‌ها و روش‌های داده‌کاوی در مقیاس بزرگ (بسیار مفید برای پژوهش و مبانی نظری).

◆ خلاصه تفاوت:

◆ " → **Kleppmann** چطور یک سیستم پایدار و مقیاس‌پذیر برای مدیریت داده طراحی کنیم؟" (تمرکز روی سیستم‌ها و معماری)

◆ " → **Leskovec et al.** چطور از داده‌های عظیم، الگو و دانش استخراج کنیم؟" (تمرکز روی الگوریتم‌ها و تحلیل داده)

Designing Data-Intensive Applications (Martin Kleppmann, 2017)

📌 تمرکز اصلی:

- ◇ معماری سیستم‌های داده‌ای مدرن
- ◇ نحوه طراحی اپلیکیشن‌هایی که باید حجم زیاد داده را با پایداری، مقیاس‌پذیری و کارایی مدیریت کنند
- ◇ 🗝️ موضوعات کلیدی:
- ◇ مقایسه **SQL vs NoSQL** و انتخاب پایگاه داده مناسب
- ◇ معماری سیستم‌های **batch processing** و **stream processing**
- ◇ replication, partitioning, consistency, CAP theorem
- ◇ طراحی سیستم‌های **distributed** و **fault-tolerant**
- ◇ 🎯 مناسب برای: مهندسان داده، معماران نرم‌افزار، کسانی که می‌خواهند **زیرساخت** بسازند یا انتخاب کنند (مثل **Kafka, Cassandra, Spark**).

Mining of Massive Datasets (Leskovec, Rajaraman, Ullman, 2020)

◆  تمرکز اصلی:

◆ الگوریتم‌ها و تکنیک‌های تحلیل و استخراج دانش از داده‌های بسیار بزرگ

◆  موضوعات کلیدی:

◆ MapReduce و مدل‌های محاسباتی داده‌های عظیم

◆ الگوریتم‌های خوشه‌بندی (clustering) و طبقه‌بندی (classification) در مقیاس بزرگ

◆ PageRank و تحلیل گراف‌های عظیم (شبکه‌های اجتماعی، وب)

◆ الگوریتم‌های data mining برای stream و big data

محتوی دور فرانس اصلی درس

1. Introduction to Big Data and MapReduce

- مدل محاسباتی MapReduce
- پیاده‌سازی الگوریتم‌ها در محیط‌های توزیع‌شده

2. Similarity Search and Locality-Sensitive Hashing (LSH)

- تعریف شباهت (Cosine, Jaccard)
- روش‌های یافتن اقلام مشابه در داده‌های عظیم

3. Data Stream Mining

- مدل داده‌های جریان (streams)
- الگوریتم‌های تقریبی برای شمارش، انتخاب نمونه، فرکانس بالا

4. Link Analysis

- الگوریتم PageRank (مبنای موتور جستجوی گوگل)
- تحلیل شبکه‌ها و گراف‌های بزرگ

5. Clustering

- الگوریتم‌های k-means, hierarchical clustering
- خوشه‌بندی در مقیاس بزرگ

6. Recommendation Systems

- فیلترینگ مشارکتی (collaborative filtering)
- مدل‌سازی ماتریسی (matrix factorization) برای سیستم‌های توصیه‌گر

7. Dimensionality Reduction

- SVD (Singular Value Decomposition)
- Principal Component Analysis (PCA)
- کاربرد در داده‌های با ابعاد بالا



8. Frequent Itemsets and Association Rules

- الگوریتم Apriori
- تحلیل سبد خرید (market-basket analysis)

9. Mining Web Advertising and Social Networks

- مدل‌های تبلیغات آنلاین
- تحلیل رفتار کاربران و شبکه‌های اجتماعی

10. Large-Scale Machine Learning Basics

- در نسخه‌های جدید، بخش‌هایی از یادگیری ماشین در مقیاس بزرگ (مثلاً regression, classification) در محیط توزیع‌شده) هم اضافه شده است.

Book: Mining of Massive Datasets

1. Foundations of Data Systems

- چگونه داده‌ها در سیستم‌های مختلف ذخیره و پردازش می‌شوند
- بررسی trade-offها بین consistency, availability, partition-tolerance (CAP theorem)
- مفاهیم data models و query languages

2. Data Storage and Retrieval

- Storage engines: Log-structured storage vs B-trees
- Indexing, replication, compaction
- پایگاه داده‌های relation و NoSQL (key-value, document, column-family)

3. Data Encoding and Evolution

- فرمت‌های داده (JSON, Avro, Protocol Buffers)
- مدیریت schema evolution در سیستم‌های توزیع شده

4. Replication

- انواع replication: single-leader, multi-leader, leaderless
- consistency models و conflict resolution

5. Partitioning (Sharding)

- تقسیم داده‌ها بین سرورها برای افزایش مقیاس‌پذیری
- طراحی partition key و مسائل مربوط به rebalancing

6. Transactions

- مفهوم ACID و انواع تراکنش‌ها در سیستم‌های توزیع شده
- Consensus protocols: Paxos, Raft

7. The Trouble with Distributed Systems

- چالش‌های سیستم‌های توزیع شده: failures, latency, network partitions
- Idempotence و retry strategies

8. Batch and Stream Processing

- پردازش داده در دسته (batch) و جریان (stream)
- مدل‌های Kappa و Lambda
- ابزارها: Hadoop, Spark, Kafka Streams

9. Derived Data and Materialized Views

- View materialization, caching, data pipelines
- ETL و event sourcing

10. Consistency and Consensus

- Strong vs weak consistency
- Distributed consensus برای هماهنگی و reliability

Book: Designing Data-Intensive Applications

ویژگی	Designing Data-Intensive Applications (Kleppmann)	Mining of Massive Datasets (Leskovec, Rajaraman, Ullman)
تمرکز اصلی	معماری و طراحی سیستم‌های داده‌ای در مقیاس بزرگ	الگوریتم‌ها و تکنیک‌های داده‌کاوی و تحلیل داده‌های عظیم
سطح	پیشرفته (معماری نرم‌افزار و مهندسی سیستم‌ها)	دانشگاهی - پیشرفته (علوم داده و الگوریتم‌ها)
موضوعات کلیدی	<ul style="list-style-type: none"> - SQL vs NoSQL - Replication, Sharding - consistency models و CAP theorem - Batch vs Stream processing - طراحی fault-tolerant systems 	<ul style="list-style-type: none"> - MapReduce و مدل محاسباتی - الگوریتم‌های خوشه‌بندی و طبقه‌بندی - PageRank و تحلیل گراف - الگوریتم‌های mining برای streams - کاهش ابعاد (dimensionality reduction)
نوع نگاه	System-oriented (چطور داده را ذخیره و پردازش کنیم در سطح زیرساخت)	Algorithm-oriented (چطور داده را تحلیل و از آن دانش استخراج کنیم)
مخاطب اصلی	مهندس داده، معمار نرم‌افزار، توسعه‌دهنده سیستم‌های توزیع‌شده	پژوهشگر، دانشجوی علوم داده/یادگیری ماشین، متخصص تحلیل داده
کاربردها	طراحی و پیاده‌سازی سیستم‌هایی مثل Kafka، Spark، Cassandra، Elasticsearch	تحلیل شبکه‌های اجتماعی، موتور جستجو، توصیه‌گرها، داده‌کاوی در وب
سبک کتاب	توضیح مفهومی + مثال‌های معماری + مقایسه ابزارها	رویکرد ریاضی/الگوریتمی + تمرینات دانشگاهی + اثبات‌ها
پیش‌نیازها	آشنایی با پایگاه داده‌ها و مفاهیم سیستم‌های توزیع‌شده	ریاضی (احتمال، آمار)، الگوریتم‌ها، برنامه‌نویسی سطح دانشگاهی
مناسب برای	کسی که می‌خواهد بداند چطور سیستم big data طراحی کند	کسی که می‌خواهد بداند چطور داده‌های عظیم را تحلیل کند

سیلابس درس

1-Basic Concepts

- Importance of Big data
- Applications of Big data
- 8v concept in Big data
- Challenges of Big data
- Tools of Big data
- Visualizations of Data
- GPU architectures in Big data
- Kappa and Lambda architectures
- MLOPS-machine learning model life cycle
- Curse of Dimensionality
- Supervised vs unsupervised learning
- TEST, TRAIN, VAL, KFOLD

2. Foundations of Data Systems

- Data types and models (relational, document, graph)
- OLTP vs OLAP systems
- consistency, availability, and durability
- Characteristics of reliable, scalable, and maintainable system

3. Unsupervised ML

- Clustering
- Clustering Problems
- select clustering algorithm for bigdata
- k-means and mean shift
- Hierarchical clustering
- Defining “nearness”
- Merging and stopping policy

4. Data Models and Query Languages

- relational vs document vs graph databases
- SQL و NoSQL
- query processing basics

5. Storage and Retrieval

- storage engines: log-structured storage, B-trees, LSM-trees
- indexing techniques
- data compression and encoding

6. Encoding and Evolution

- serialization formats (JSON, Avro, Protocol Buffers)
- schema evolution
- backward/forward compatibility

7. The Trouble with Distributed Systems

- clock synchronization and causality
- fault tolerance, network partitions, timeouts
- consistency vs availability trade-offs (CAP theorem)

8. Batch and Stream Processing

- batch processing (Hadoop, MapReduce)
- stream processing (Spark Streaming, Kafka Streams)
- exactly-once vs at-least-once semantics

9. Derived Data

- materialized views
- caches, indexes, search engines
- change-data-capture and event sourcing

10. Similarity Search and Locality-Sensitive Hashing (LSH)

- Near-duplicate detection (e.g., web pages, documents)
- Minhashing and shingling techniques
- LSH for approximate nearest neighbors

11. Data Stream Mining

- Stream models and challenges
- Sampling techniques
- Counting distinct elements, frequency moments, heavy hitters

12. Link Analysis

- Graphs and networks basics
- PageRank algorithm
- HITS and other ranking methods

13. Dimensionality Reduction

- SVD (Singular Value Decomposition)
- Principal Component Analysis (PCA)
- Graphical model (MPPCA, MFA)

14. Supervised Learning at Scale

- Decision trees, SVMs, and logistic regression
- Stochastic gradient descent (SGD)
- Handling imbalanced datasets

15. Partitioning (Sharding)

- partitioning strategies: range, hash, composite
- rebalancing shards
- handling hotspots

16. Deep learning at scale

- Serving a TensorFlow Model, TFLite, Mobile or Embedded Device
- Serving through the gRPC and REST API
- Scaling up TF Serving and Parallelized execution
- Using GPUs to Speed Up Computations- Managing the GPU RAM
- Training Models Across Multiple Devices (Model Parallelism vs Data parallelism)
- TensorFlow Cluster
- Transformers
- Mamba |

بارم بندی

بارم بندی درس (۲۴ نمره)

◆ نمرات اصلی

◆ پایان ترم ۹ نمره

◆ تکالیف ۷ نمره

◆ پروژه پایانی ۲ نمره

◆ ارایه مقاله ۲ نمره

◆ نمرات امتیازی

◆ تکالیف اختیاری ۲ نمره

◆ ارزیابی استاد از دانشجو ۱ نمره اضافی (به صورت استثنا)

◆ تولید محتوی ۱ نمره اضافی

موضوعات تکالیف

- ۱- پیاده سازی برخی از الگوریتم ها با زبان برنامه نویسی Python
- ۲- کار با بسترهایی مثل Kafka
- ۳- کار با بسترهایی مثل Spark و Hadoop
- ۴- پیاده سازی یک الگوریتم بر روی بستر GPU با استفاده از CUDA
- ۵- انجام Fine tuning یک شبکه LLM متن باز
- ۶- اجرای یک کد با استفاده از شبکه Mamba