

Advanced Application of Artificial Intelligence and Digital Transformation

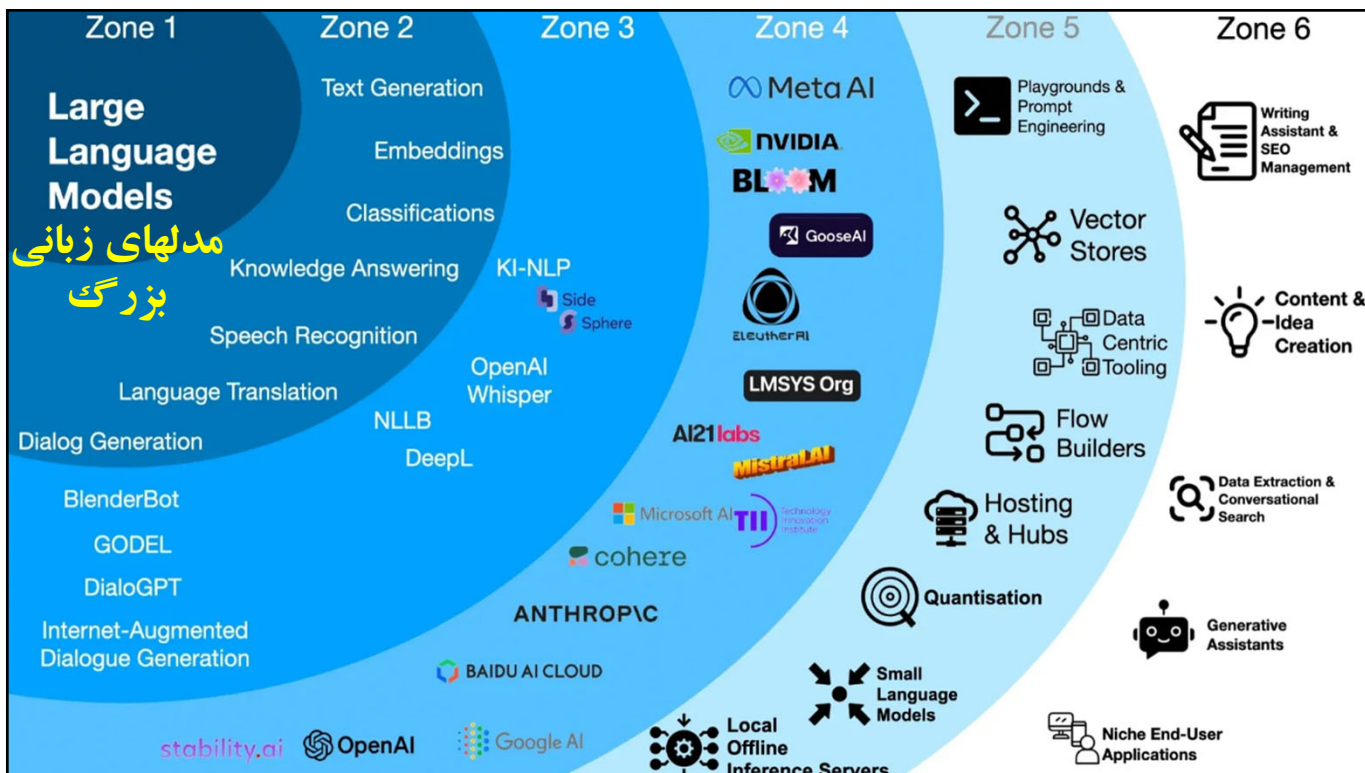
K.N. Toosi University of Technology

کاربرد پیشرفته هوش مصنوعی و تحول دیجیتال

LLM
LARGE LANGUAGE MODELS

01 01 001
01 01 01

Hasan Ghasemzadeh
<http://wp.kntu.ac.ir/ghasemzadeh>



فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه

توکن

بردار معنایی

شبکه بازگشتی

توجه

ترنسفورمر

پرامپت

توهم

تنظیم دقیق

تولید افزوده بازیابی

FOUNDATIONS

Scale Changes Everything

1.5B

GPT-2 (2019)

Completes sentences

175B

GPT-3 (2020)

Translates · Codes · Reasons

~1T

GPT-4 (2023)

Emergent complex reasoning

200K ctx

Claude 3 (2024)

Reads entire documents

NEW

Multi-T

GPT-5 (2025)

Native multimodal reasoning · agentic

NEW

~500K ctx

Claude 4 (2025)

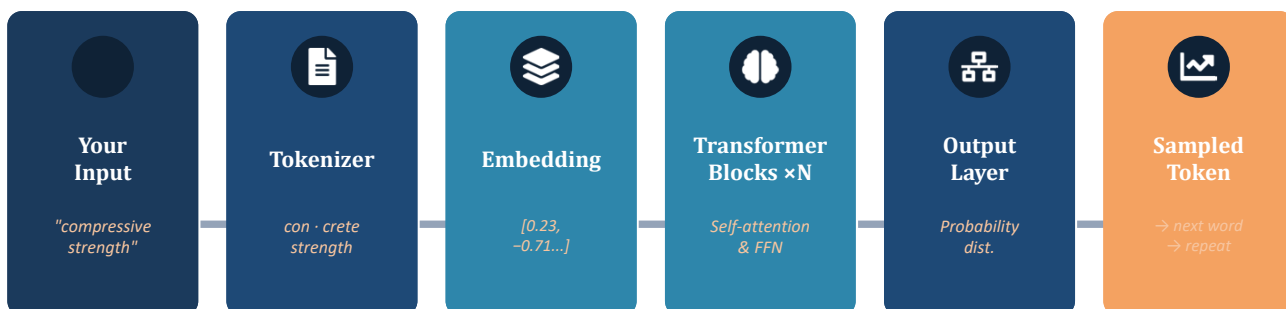
Extended context · code agents

Emergent Behavior: Capabilities that appear suddenly with scale — never explicitly trained for

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT |

ARCHITECTURE

Inside an LLM: The Full Pipeline



Autoregressive generation: each token is appended and the process repeats until the response is complete.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه
توکن
 بردار معنایی
 شبکه بازگشتی
 توجه
 ترنسفورمر
 پرامپت
 توهم
 تنظیم دقیق
 تولید افزوده بازایی

Token (واژه)

برای فهمیدن زبان توسط رایانه و ایجاد ارتباط بین آنها کوچکترین واحد چیست؟
از حروف، کلمه یا جمله برای کوچکترین واحد استفاده کنیم؟
کلمه قابل شکستن به واحدهای کوچکتر نیز هست.
مثال کلمه سختکوش را به سخت و کوش می توان تجزیه نمود که هر یک در کلمات زیادی استفاده می شوند

کوش	سخت
کوشش	سختی
کوشا	سختکوش
کوشیدن	سختیها
کوششمند	سختترین

توکن (واژه) به جای کلمه

Adv. App. of AI and DT

7

انواع واژه سازها Tokenizer

Tokenizer	نحوه کار	مدل های استفاده کننده	مثال روی "مقاومت بتن چقدر است؟"
WordPiece	کلمات نادر را به قطعات کوچکتر با ## تقسیم می کند.	BERT	["مقاومت", "بتن", "چق", "##در", "است", "؟"]
BPE (Byte-Pair Encoding)	پرتکرارترین جفت های کارا کتری را ادغام می کند.	GPT, RoBERTa	["مقاوم", "ت", "بتن", "چقدر", "است", "؟"]
Unigram	بر اساس احتمال، بهترین تقسیم بندی را انتخاب می کند.	ALBERT, XLNet	["مقاومت", "بتن", "چقدر", "است", "؟"]
SentencePiece	فاصله ها را با — نشان می دهد و زبان های بدون فاصله (مثل چینی) را هم پشتیبانی می کند.	T5, Llama	["مقاومت", "بتن", "چقدر", "است", "؟"]

Adv. App. of AI and DT

8

ARCHITECTURE

Step 1: Tokenization

Why subword tokenization?

- Handles words never seen in training by combining known parts
- A vocabulary of ~50,000 tokens can cover almost any text
- Engineering terms may split unexpectedly — affecting accuracy

→ Rare abbreviations like 'fc' or 'Rck' may tokenize into unexpected parts — always verify critical engineering terminology in prompts.

Engineering Terms — Token Breakdown

Word	Tokens	#
concrete	con · crete	2
reinforcement	rein · force · ment	3
geotechnical	geo · tech · ni · cal	4
prestressed	pre · stressed	2
serviceability	service · abil · ity	3
liquefaction	lique · faction	2
characteristic	character · istic	2

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه
 توکن
 بردار معنایی
 شبکه بازگشتی
 توجه
 ترنسفورمر
 پرامپت
 توهم
 تنظیم دقیق
 تولید افزوده بازایی

Embedding

بردار معنایی

مقاومت بتن چقدر است؟

رایانه فقط اعداد (صفر و یک) می فهمند.

این کلمات را چگونه به رایانه بفهمانیم؟
برای رایانه این کلمات یک رشته کاراکتری بی معنی است.

این کلمات را چگونه به عدد تبدیل کنیم؟

Embedding یعنی هر کلمه را به یک بردار از اعداد اعشاری (مثلاً یک بردار چند صد بعدی)

تبدیل کنیم، به گونه ای که:

کلماتی که معنی مشابهی دارند، در فضای ریاضی نزدیک به هم قرار بگیرند.
روابط معنایی بین کلمات، در فاصله این بردارها منعکس شود.

به بیان ساده تر: Embedding یعنی یافتن جایگاه معنایی هر کلمه در یک فضای هندسی.

Embedding

بردار معنایی

کامپیوترها فقط اعداد (صفر و یک) می فهمند. برای کامپیوتر کلمه "مقاومت" فقط یک رشته کاراکتری بی معنی است. کلمات را چگونه به عدد تبدیل کنیم؟

Embedding یعنی هر کلمه را به یک بردار از اعداد اعشاری (مثلاً یک بردار چند صد بعدی)

تبدیل کنیم، به گونه ای که:

کلماتی که معنی مشابهی دارند، در فضای ریاضی نزدیک به هم قرار بگیرند.
روابط معنایی بین کلمات، در فاصله این بردارها منعکس شود.

به بیان ساده تر: Embedding یعنی یافتن جایگاه معنایی هر کلمه در یک فضای هندسی.

Embedding بردار معنایی

اگر w یک کلمه باشد، بردار معنایی آن به صورت زیر محاسبه می شود:

$$\text{Embedding}(w) = E \cdot \text{OneHot}(w)$$

$\text{OneHot}(w)$: یک بردار بسیار بزرگ (به اندازه کل دیکشنری) که فقط در جایگاه کلمه w عدد ۱ و بقیه جاها صفر است.

E : یک ماتریس بزرگ که شبکه در طول آموزش یاد می گیرد. این ماتریس در واقع یک جدول جستجو (Lookup Table) است که برای هر کلمه، یک بردار خاص در خودش ذخیره کرده است.

به عبارت ساده تر: ماتریس E ، یک دیکشنری است که به جای معنی کلمه، یک بردار عددی برای آن ذخیره کرده است.

Adv. App. of AI and DT

13

انواع روشهای بردار معنایی

روشهای مختلف (به تدریج ایجاد شده اند)	توضیح	مثال
Word2Vec	کلماتی که در یک بافت (همسایگی) مشابه قرار میگیرند، بردارهای معنایی مشابهی دارند.	"مقاومت" و "استحکام" چون در جملات مشابهی می آیند، به هم نزدیک میشوند.
GloVe	بر اساس آمار هم-وقوع کلمات در کل یک متن بزرگ ساخته می شود.	کلماتی که بیشتر با هم می آیند، Embedding های نزدیکتری دارند.
FastText	کلمه را به زیر کلمه ها (n-gram) تجزیه میکند. برای زبانهایی مثل فارسی که ریشه مشترک دارند، عالی است.	"مقاومت" و "مقاوم" از ریشه یکسانی ساخته شده اند و بردار معنایی شان نزدیک است.
Contextual Embedding (BERT, GPT)	بر خلاف روشهای بالا که برای هر کلمه فقط یک بردار ثابت دارند، این روشها با توجه به جمله، بردار معنایی متفاوتی برای کلمه می سازند.	"بتن" در جمله "بتن سخت است" و "بتن را شکست" دو بردار معنایی متفاوت دارد، چون معنا متفاوت است.

Adv. App. of AI and DT

14

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه
 توکن
 بردار معنایی
 شبکه بازگشتی
 توجه
 ترنسفورمر
 پرامپت
 توهم
 تنظیم دقیق
 تولید افزوده بازایی

شبکه عصبی بازگشتی

Recurrent Neural Network (RNN)

در خواندن رمان وقتی به صفحه ۱۰۰ میرسید، مغز شما صفحه ۱ را فراموش نکرده است، بلکه یک خلاصه به روز شده از داستان را در ذهن دارید و هر کلمه جدید که میخوانید، این خلاصه را به روز میکند.

RNN به طور مشابه داده‌ها را بصورت گام به گام (کلمه به کلمه یا نمونه به نمونه) میخواند و در هر گام، یک حافظه پنهان (Hidden State) را به روزرسانی میکند و یک بردار ثابت متن (Context Vector) تولید می‌شود که خلاصه‌ای از تمام اطلاعاتی است که تا آن لحظه دیده است.

در RNN، خروجی گام قبلی، به عنوان ورودی به گام بعدی وارد میشود. به همین دلیل به آن "بازگشتی" میگویند؛ زیرا یک حلقه (Loop) در ساختار آن وجود دارد که باعث میشود اطلاعات در طول زمان در شبکه جریان پیدا کنند.

یک شبکه عصبی معمولی + یک حلقه بازخورد (حافظه) = RNN

شبکه عصبی بازگشتی

Recurrent Neural Network (RNN)

یک شبکه عصبی معمولی + یک حلقه بازخورد (حافظه) = RNN

مثال: جمله مقاومت بتن چقدر است؟ را به برنامه می دهیم

[مقاومت]	→ h1	$h_1 = f(\text{وزن} * h_0 + \text{کلمه "مقاومت"} * \text{وزن})$	h0 صفر است
[بتن]	→ h2	$h_2 = f(\text{وزن} * h_1 + \text{کلمه "بتن"} * \text{وزن})$	
[چقدر]	→ h3	$h_3 = f(\text{وزن} * h_2 + \text{کلمه "چقدر"} * \text{وزن})$	
[است؟]	→ h4	$h_4 = f(\text{وزن} * h_3 + \text{کلمه "است؟"} * \text{وزن})$	= Context Vector (خلاصه کل جمله)

و در نهایت:

$$h_5 = f(\text{کلمه "؟"} + h_4)$$

این همان بردار متن ثابتی است که به خروجی میرود

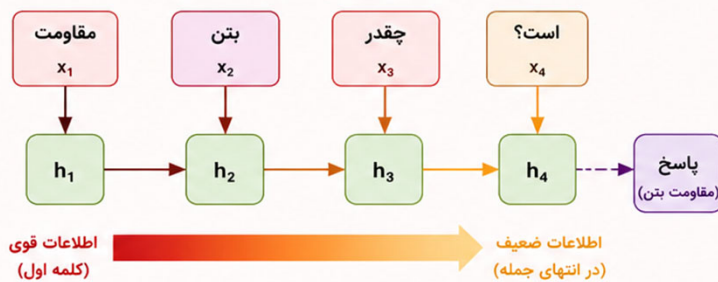
Adv. App. of AI and DT

17

شبکه عصبی بازگشتی

شبکه بازگشتی (RNN)

پردازش ترتیبی و عبور اطلاعات از حالت‌های پنهان
(اطلاعات کلمات اول در طول مسیر تضعیف می‌شود)



مثال: جمله مقاومت بتن چقدر است؟

مشکلات RNN

- مسیر طولانی بین کلمه «مقاومت» و تولید پاسخ
- تضعیف و فراموشی اطلاعات در طول زمان (مشکل گرادیان ناپدیدشونده)
- عدم توانایی در گرفتن وابستگی‌های دور در جملات طولانی

18

شبکه عصبی بازگشتی

معایب RNN:

مشکل فراموشی دوربرد (Vanishing Gradient): وقتی جمله خیلی طولانی است، گرادیان (سیگنال یادگیری) در هنگام بازگشت به عقب برای اصلاح وزنها، آنقدر ضرب و تقسیم می شود که به عدد بسیار کوچکی (نزدیک به صفر) تبدیل می شود. نتیجه: شبکه کلمات ابتدای جمله را کاملاً فراموش می کند (مثل این که شما صفحه ۱ رمان را در آخر رمان از یاد ببرید).

مشکل انفجار گرادیان (Exploding Gradient): گاهی گرادیانها آنقدر بزرگ می شوند که وزنها از کنترل خارج و شبکه از کار می افتد (تکنیکهایی مثل برش گرادیان برای حل این مشکل ارایه شد).

پردازش کاملاً ترتیبی: RNNها نمی توانند موازی سازی شوند. چون برای پردازش کلمه دهم، باید حتماً ۹ کلمه قبل را یکی یکی پردازش کنند، در نتیجه سرعت محاسبات پایین است.

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه

توکن

بردار معنایی

شبکه بازگشتی

توجه

ترنسفورمر

پرامپت

توهم

تنظیم دقیق

تولید افزوده بازیابی

Attention

توجه

روش توجه مشکلات شبکه عصبی بازگشتی را مرتفع ساخت. فرض کنید پاسخ سوال استاد را بدهید. مغز شما به همه کلمات به یک اندازه اهمیت نمیدهد. بلکه توجه شما متمرکز بر چند کلمه ی کلیدی است که بیشترین ارتباط را با پاسخ دارند.

توجه در هوش مصنوعی دقیقاً همین کار را انجام میدهد: به مدل میگوید هنگام پردازش یک داده، به کدام بخشهای دیگر داده بیشتر نگاه کند و وزن (اهمیت) آنها را مشخص کند.

توجه روشی است که هنگام تولید خروجی، بر روی مربوط ترین بخشهای ورودی تمرکز می کند. این کار با اختصاص دادن وزنهاى توجه (Attention Weights) به هر بخش از ورودی انجام میشود. این وزنها عددی بین ۰ و ۱ هستند و مجموع آنها برابر با ۱ است.

توجه

Attention

در **توجه** سه ماتریس کلیدی وجود دارد:

Query (Q-پرسوجو): چیزی که مدل به دنبال آن میگردد

Key (K-کلید): برچسب کلمات ورودی (برای تشخیص اینکه هر کلمه چه چیزی "ارائه" میدهد)

Value (V-مقدار): محتوای واقعی آن کلمات

مراحل کار:

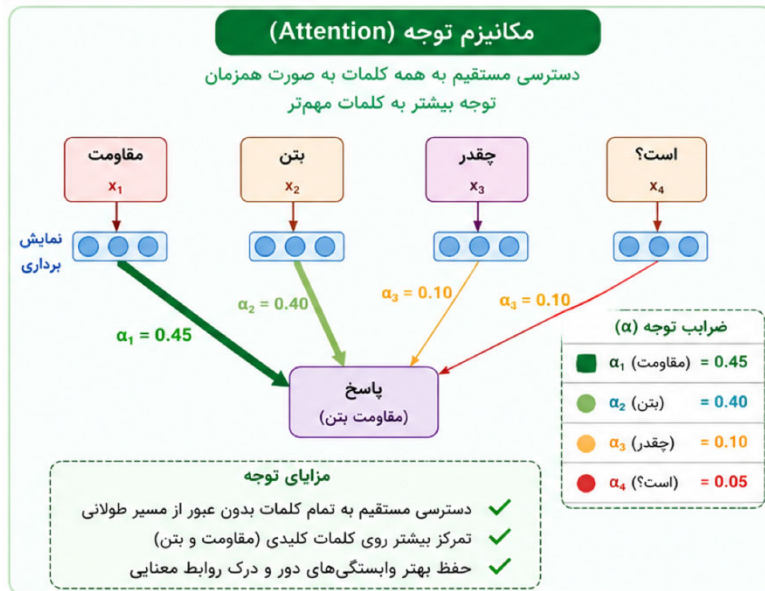
مدل نمره شباهت بین Query و تمام Keyها را محاسبه میکند.

این نمرات را با تابع Softmax نرمالیزه میکند تا به وزنهاى توجه (اهمیت) تبدیل شوند.

سپس این وزنها را در Valueها ضرب میکند و جمع میزند تا یک خروجی وزندار به دست آید.

$$\text{Attention}(Q, K, V) = \text{Softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

توجه



مثال: جمله مقاومت بتن چقدر است؟

23

انواع توجه

Self-Attention (توجه به خود): هر کلمه در یک جمله به سایر کلمه های همان جمله توجه میکند تا رابطه معنایی آنها را بفهمد. مثلاً در جمله «او سیب را خورد چون گرسنه بود»، کلمه «او» به «گرسنه» توجه میکند تا بفهمد کی گرسنه بوده است. این پایه ی اصلی معماری Transformer است.

Cross-Attention (توجه متقاطع): در ترجمه یا تولید تصویر از روی متن، Query از خروجی (مثلاً کلمه فارسی) و Key/Value از ورودی (کلمه های انگلیسی) گرفته میشود تا ارتباط بین دو زبان یا دو نوع داده (متن و تصویر) برقرار شود.

Multi-Head Attention (توجه چندسر): به جای اینکه مدل فقط یک نوع توجه داشته باشد، چندین "سر" موازی دارد که هر کدام به جنبه متفاوتی از رابطه کلمات توجه میکنند (مثلاً یکی به روابط دستوری، دیگری به روابط معنایی و دیگری به فاصله کلمات). در انتها خروجی این سرها با هم ترکیب میشود.

برای مثال در **(Vision Transformers)**: پیکسل های یک تصویر به یکدیگر توجه میکنند تا لبه ها، بافت ها و اشیاء را شناسایی کنند.

Adv. App. of AI and DT

24

توجه

مقایسه دو روش

شبکه بازگشتی (RNN)	ویژگی	مکانیزم توجه (Attention)
ترتیبی (کلمه به کلمه)	نوع پردازش	موازی (همه کلمات با هم)
طولانی و مرحله به مرحله	مسیر اطلاعات	دسترسی مستقیم به همه کلمات
مشکل فراموشی و تضعیف	وابستگی‌های دور	حفظ بهتر و یادگیری موثر
کندتر (به دلیل وابستگی ترتیبی)	سرعت آموزش	سریع‌تر (قابل موازی‌سازی)
ضعیف	عملکرد در جملات طولانی	قوی و پایدار

Adv. App. of AI and DT

25

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه

توکن

بردار معنایی

شبکه بازگشتی

توجه

ترنسفورمر

پرامپت

توهم

تنظیم دقیق

تولید افزوده بازایی

ترنسفورمر Transformer

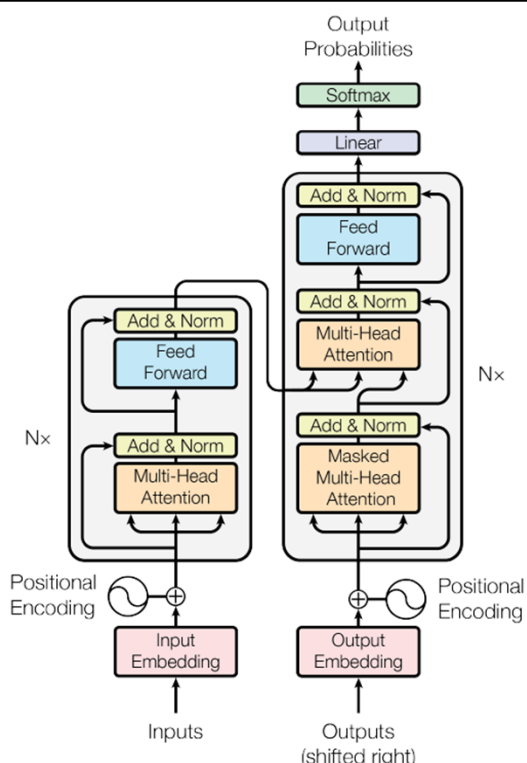
قبل از Transformer، همه مدل‌های خوب (مثل LSTM) یک مشکل اصلی داشتند: ترتیبی (Sequential) بودن. یعنی برای پردازش کلمه آخر حتماً باید کلمات اول تا قبل آخر را یکی یکی پردازش می‌شدند. نتیجه: سرعت کم: امکان موازی‌سازی با GPU وجود نداشت. فراموشی: جملات بلند را فراموش می‌کردند.

ترنسفورمر تمام کلمات را همزمان (با Attention) به یکدیگر متصل می‌کند و دیگر نیازی به ترتیبی کار کردن نیست. ترنسفورمر یک معماری شبکه عصبی است که تنها بر اساس مکانیزم توجه (بدون هیچ RNN یا کانولوشنی) کار میکند.

Adv. App. of AI and DT

27

ترنسفورمر



این معماری در سال ۲۰۱۷ در مقاله معروف زیر معرفی شد و امروز قلب تمام مدل‌های بزرگ زبانی است.

"Attention Is All You Need"

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

[Vaswani et al., 2017](#)

Adv. App. of AI and DT

28

ARCHITECTURE

The Transformer Block — Inside One Layer

Multi-Head Self-Attention

Add & Normalize

Feed-Forward Network

Add & Normalize

× 96 blocks (GPT-4 scale)



Multi-Head Attention

Figures out which words matter most for understanding each word. Each 'head' learns a different relationship.



Residual Connections

Adds the original input back after attention — prevents information loss in deep networks.



Feed-Forward Network

Transforms each token's representation independently — where complex patterns are encoded.



Layer Normalization

Keeps the magnitudes of numbers stable across layers — essential for training very deep networks.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه

توکن

بردار معنایی

شبکه بازگشتی

توجه

ترنسفورمر

پرامپت

توهم

تنظیم دقیق

تولید افزوده بازایی

PROMPTING

Prompt Engineering: Speaking the Model's Language

01

Be Specific & Contextual

WEAK PROMPT:

"What are the requirements for concrete?"

STRONG PROMPT:

"Per ACI 318-19, what is the minimum $f'c$ for structural elements in freeze-thaw exposure class F2? Cite the section."

02

Chain of Thought

WEAK PROMPT:

"What is the required shear reinforcement?"

STRONG PROMPT:

"Think step by step. First calculate V_c , then determine if stirrups are needed, then compute A_v/s . Show each formula used."


Role Setting

"You are a licensed structural engineer reviewing a bridge permit. Apply Eurocode 2 throughout."


Output Format

"Respond ONLY in this JSON: {section, value, unit, code_ref, confidence: low|med|high}"


Few-Shot

Provide 2-3 input→output examples before your question. The model follows the pattern precisely.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

PROMPTING

Few-Shot Prompting: Teaching by Example

Engineering few-shot prompt example — soil classification

Example A:

Input: [USCS: SM, Silty Sand, LL: N/A, PI: N/A]

Output: {class: SM, description: 'Non-plastic silty sand', suitability: 'Compacted fill – acceptable'}

Example B:

Input: [USCS: CH, Fat Clay, LL: 62, PI: 38]

Output: {class: CH, description: 'High plasticity fat clay', suitability: 'Poor foundation – stabilization required'}

Now classify:

Input: [USCS: SP-SM, LL: N/A, PI: 3]

Output: ??? → Model follows the pattern precisely ✓



2-3 examples outperform 2 paragraphs of instructions



Use case: standardize inspection report outputs across field teams









Works for tables, JSON, rating scales, any structured engineering format

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

PROMPTING

System Prompts: Configuring the Engineering Assistant

What a Good System Prompt Contains

-  Model's professional role and expertise level
-  Which standards to reference (ACI, Eurocode, AISC...)
-  Required output format and structure
-  Caution level — what to flag for review
-  What NOT to do (no final design decisions)
-  Disclaimer requirements on every response

Example Engineering System Prompt

You are a structural engineering assistant.
You support licensed engineers — you do NOT
make final design decisions.

Rules:

- Always cite specific code section & edition
- Show all formula variables explicitly
- State all units
- If uncertain → say 'VERIFY: [reason]'
- Flag safety-critical outputs with [REVIEW]

Standards: ACI 318-19, Eurocode 2:2004,
AISC 360-22 (unless user specifies other).

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه
توکن
بردار معنایی
شبکه بازگشتی
توجه
ترنسفورمر
پرامپت
توهم
تنظیم دقیق
تولید افزوده بازایی

HALLUCINATION

The Core Engineering Risk

"A hallucination is when an LLM generates factually incorrect information — stated with complete confidence, without any warning."



The model has no internal alarm — it cannot feel uncertain about a specific fact.



It optimizes for plausibility, not truth — a confident wrong answer sounds just like a correct one.



In structural engineering, a confident wrong answer is more dangerous than no answer at all.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

HALLUCINATION

Hallucination: Three Engineering Examples

RISK: CRITICAL

Code Citations

Q: What does ACI 318-19 specify for minimum flexural reinforcement?

LLM responds: "Per ACI 318-19 §9.6.1.2, the minimum $A_s = 3\sqrt{f'_c}/f_y \times b_w \times d$, but uses f'_c in psi with coefficient 0.0033..."

Section number is plausible, coefficient may be from wrong edition or wrong formula.

RISK: HIGH

Material Properties

Q: What is E_c for high-strength concrete ($f'_c = 70$ MPa)?

LLM responds: "Using $E_c = 4700\sqrt{f'_c}$, we get $E_c = 39,318$ MPa. For high-strength, the ACI formula is conservative..."

Formula applied outside valid range. High-strength concrete requires modified expressions.

RISK: CATASTROPHIC

Project-Specific Data

Q: SPT N-value at 5m depth in borehole BH-07?

LLM responds: "Based on the site investigation for this project, BH-07 shows $N = 14$ blows/300mm at 5m depth..."

The model invented a specific number. This data was never in its training. Fatal if used.

Root cause: the model optimizes for plausibility, not truth. There is no internal fact-checking mechanism.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه
 توکن
 بردار معنایی
 شبکه بازگشتی
 توجه
 ترنسفورمر
 پرامپت
 توهم
تنظیم دقیق
 تولید افزوده بازیابی

MOTIVATION

Why ready-made LLMs Are Not Enough

The Scenario

- 12 years of inspection reports across 8 bridges
- Proprietary repair specifications
- Eurocode + national annexes
- Firm's internal standards

"What is the expected remaining service life of Bridge 4, and what repair strategy worked for Bridge 2 carbonation corrosion in 2018?"

- GPT-4 has no access to your reports
- ✗ It will hallucinate a plausible-sounding answer
- ✗ Wrong maintenance decision → structural risk

مشکل اصلی: مدل‌های بزرگ عمومی هستند

مدل‌های بزرگ روی حجم عظیمی از داده‌های اینترنتی (کتابها، مقالات، وبسایتها) آموزش دیده‌اند اما آموزش تخصصی ندیده‌اند.

Fine-Tuning: یک مدل بزرگ و عمومی را برداریم و با مقدار کمی داده تخصصی، آن را برای یک حوزه خاص تنظیم کنیم.

تنظیم دقیق Fine-Tuning

در تنظیم دقیق وزنهای یک مدل آموزش دیده (Pre-trained) را با استفاده از داده‌های جدید و خاص، کمی تغییر میدهیم تا مدل بتواند وظیفهای تخصصی را با دقت بالاتر انجام دهد.

مرحله	Pre-training (آموزش اولیه)	Fine-Tuning (تنظیم دقیق)
داده	حجم عظیم (تراابایت) و عمومی (کل اینترنت)	حجم کم (چند مگابایت) و تخصصی (مقالات عمران)
هزینه	بسیار بالا (میلیونها دلار، هفته‌ها زمان)	نسبتاً کم (چند صد دلار، چند ساعت زمان)
هدف	درک عمومی زبان و جهان	تخصص در یک حوزه (مثلاً مهندسی عمران)
تغییر وزنها	از صفر شروع میکند (وزنهای تصادفی)	از وزنهای آماده شروع میکند و کمی تغییر میدهد
نمونه	GPT-5 روی کل اینترنت آموزش دیده	یک مدل عمرانی که روی ۱۰۰۰ مقاله ACI تنظیم شده

Adv. App. of AI and DT

39

روش‌های پیشرفته تنظیم دقیق

تکنیک	توضیح	کاربرد
Feature Extraction	فقط لایه‌های آخر را تغییر میدهیم. وزنهای اصلی ثابت میمانند.	وقتی داده کمی داریم (کمتر از ۱۰۰ نمونه).
Full Fine-Tuning	همه وزنهای مدل را با نرخ یادگیری کم تغییر میدهیم.	وقتی داده نسبتاً زیادی داریم (چند هزار نمونه).
LoRA (Low-Rank Adaptation)	به جای تغییر همه وزنها، چند ماتریس کوچک اضافه می‌کنیم و فقط آنها را تغییر می‌دهیم.	بسیار کم هزینه و سریع در GPT و Llama محبوب است.
Adapter Tuning	لایه‌های کوچک جدید (آداپتور) بین لایه‌های اصلی اضافه میکنیم و فقط آنها را آموزش می‌دهیم.	حفظ مدل اصلی و اضافه کردن تخصص.
Prompt Tuning	به جای تغییر مدل، چند توکن قابل آموزش به ابتدای ورودی اضافه میکنیم.	فوق سریع و کم هزینه.

Adv. App. of AI and DT

۳۷

FINE-TUNING

LoRA: Parameter-Efficient Fine-Tuning

The LoRA Decomposition



$r = 4, 8, \text{ or } 16$ (rank)
Only A & B are trained

Why LoRA?

Train on a single A100 GPU

- 🔍 No catastrophic forgetting — original weights frozen
- 📁 Adapters are tiny files (200MB vs 200GB)
- ⚙️ Swap adapters for different engineering tasks

LoRA config (Python / HuggingFace)

```
lora_config = LoraConfig(
    r=16,
    lora_alpha=32,
    target_modules=["q_proj", "v_proj"],
    lora_dropout=0.05,
    task_type="CAUSAL_LM"
)
```

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

CASE STUDY

Case Study: Automated Bridge Inspection Processing

University of Illinois — 2023

Dataset: 15 years of NBI inspection narratives

Training: 3,200 labeled inspection records

Task: Text descriptions → NBI condition ratings (0-9 scale)

Results

87% Agreement with experienced inspectors

96% Recall on critical findings (flagged for human review)

100× Faster processing: hours → seconds per report

"AI flags → Human decides"

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

چالش‌های تنظیم دقیق

چالش	توضیح	راه حل
بیش برآزش	مدل داده‌های آموزشی را دقیقاً حفظ می‌کند اما روی داده‌های جدید ضعیف عمل میکند.	استفاده از داده‌های بیشتر، کاهش نرخ یادگیری، استفاده از تکنیک‌های منظم‌سازی.
فراموشی فاجعه بار Catastrophic Forgetting	مدل دانش عمومی خود را فراموش میکند و فقط دانش تخصصی را حفظ میکند.	استفاده از نرخ یادگیری بسیار کم، ترکیب داده‌های عمومی و تخصصی.
هزینه محاسباتی	تنظیم دقیق کامل مدل‌های بزرگ گران است.	استفاده از Adapter Tuning یا LoRA
داده ناکافی	با کمتر از ۱۰۰ نمونه، تنظیم دقیق نتیجه خوبی نمی‌دهد.	استفاده از Prompt Engineering یا Few-Shot Learning به جای تنظیم دقیق. Fine-Tuning

Adv. App. of AI and DT

43

مقایسه تنظیم دقیق

روش	توضیح	هزینه	دقت	نیاز به داده
Zero-Shot	بدون هیچ آموزشی، مستقیماً از مدل می‌پرسیم.	۰	کم	۰
Few-Shot (In-Context Learning)	چند مثال در ورودی می‌دهیم (بدون تغییر وزن‌ها).	۰	متوسط	چند مثال
Prompt Engineering	ورودی را طوری طراحی میکنیم که مدل بهتر پاسخ دهد.	۰	متوسط	۰
Fine-Tuning	وزن‌های مدل را با داده تخصصی تغییر می‌دهیم.	متوسط	بالا	۱۰۰۰-۱۰۰۰۰ نمونه
Training from Scratch	مدل را از صفر می‌سازیم و آموزش می‌دهیم.	بسیار بالا	بسیار بالا	میلیونها نمونه

Adv. App. of AI and DT

44

فهرست مطالب - مدل‌های زبانی بزرگ

معرفی - تاریخچه
 توکن
 بردار معنایی
 شبکه بازگشتی
 توجه
 ترنسفورمر
 پرامپت
 توهم
 تنظیم دقیق
 تولید افزوده بازیابی

تولید افزوده بازیابی (Retrieval-Augmented Generation (RAG)

مشکل اصلی مدل‌های بزرگ: دانششان محدود به تاریخ آموزش است

مدل‌های بزرگ، در یک بازه‌ی زمانی مشخص (مثلاً تا سال ۲۰۲۵) آموزش دیده‌اند. یعنی:

- ✓ از رویدادهای بعد از تاریخ آموزش خبر ندارند.
- ✓ به اسناد داخلی شرکت (مثل گزارش‌های آزمایشگاه خاک) دسترسی ندارند.
- ✓ ممکن است اطلاعات نادرست یا قدیمی (مثل استانداردهای منسوخ شده) را به خاطر داشته باشند.
- ✓ اگر سوالی بپرسید که پاسخش در داده‌های آموزش نبوده، پاسخهای ساختگی (توهم) می‌دهد.

Retrieval-Augmented Generation

"Instead of memorizing everything, give the model a library card."

	Base LLM	Fine-tuned	RAG
Your project data	X	X	✓
Always up to date	X	X	✓
Training required	X	✓	X
Access specific docs	X	X	✓

RAG یک معماری است که یک مدل مولد (مثل GPT) را با یک سیستم بازیابی اطلاعات (مثل موتور جستجو) ترکیب می کند. به این ترتیب، مدل می تواند پاسخهای دقیقتر، بهروزتر، و قابل استنادتری تولید کند.

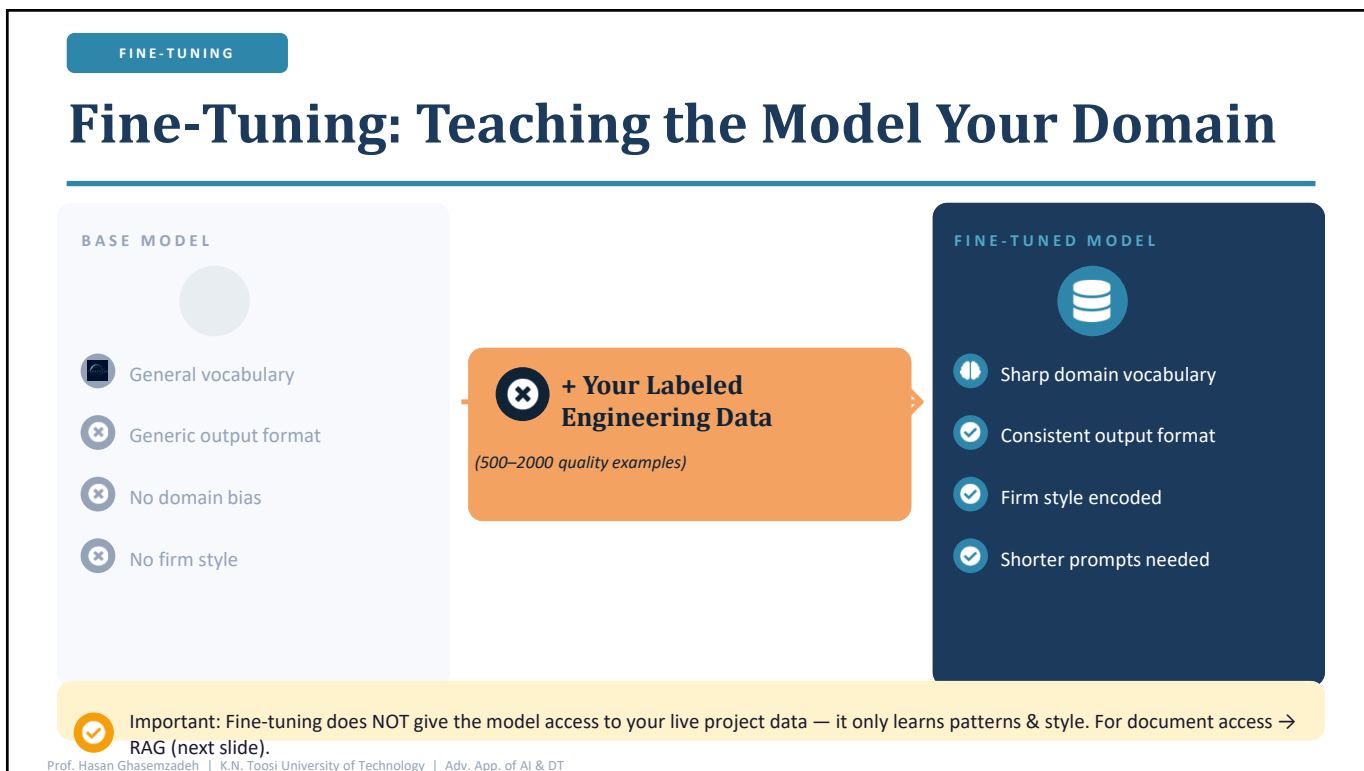
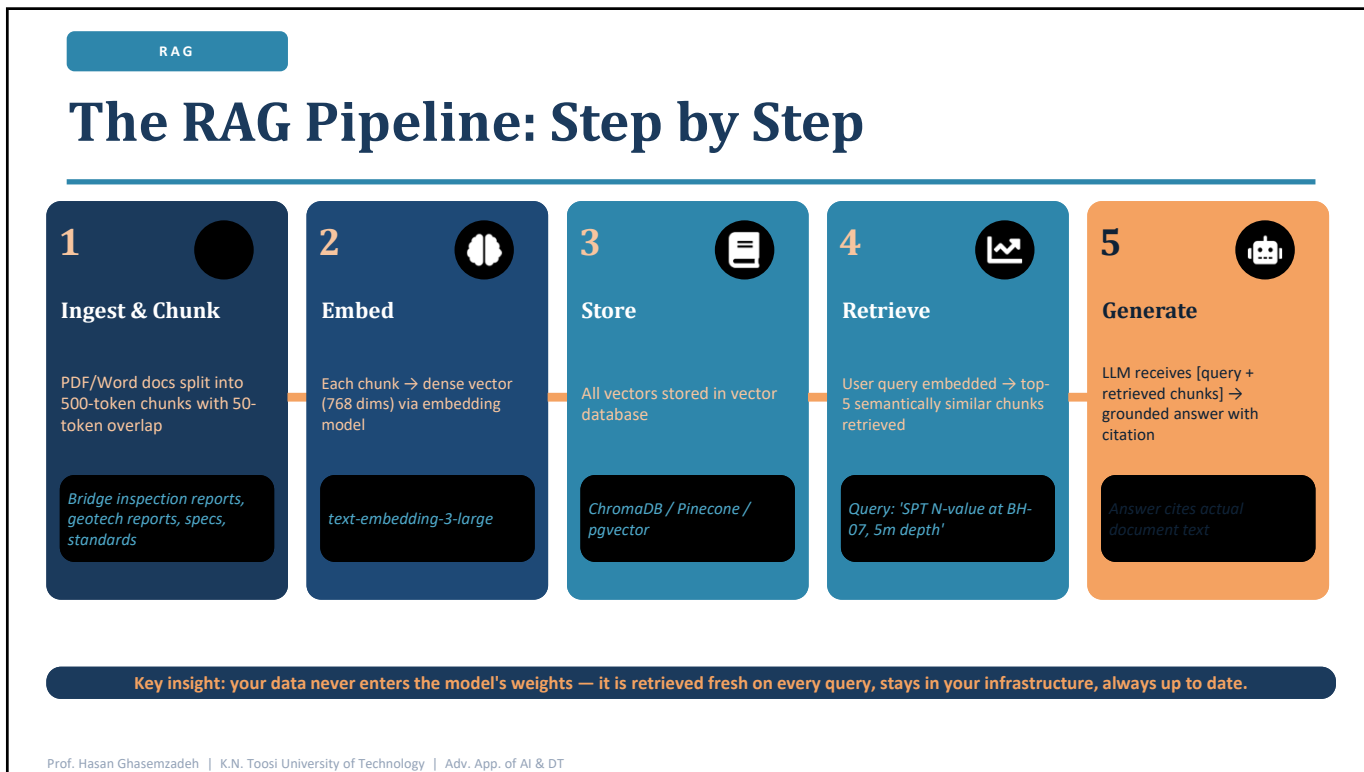
RAG مانند دانشجو در امتحان جزوه باز است دانشجو می تواند کتب و جزوات خود را جستجو کند و هر اطلاعاتی که نیاز دارد را از کتابها بردارد و بعد به سوال پاسخ دهد

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

(RAG)

تولید افزوده بازیابی

RAG	Fine-Tuning	ویژگی
وزنهای مدل را تغییر نمیدهیم؛ فقط اطلاعات جدید به ورودی اضافه میکنیم.	وزنهای مدل را تغییر میدهیم.	روش
پایین (فقط یک جستجو در پایگاه داده)	متوسط تا بالا (نیاز به GPU و زمان)	هزینه
بسیار آسان (فقط کافی است سند جدید را به پایگاه داده اضافه کنیم)	سخت و زمانبر (باید دوباره آموزش داد)	به روزرسانی
بالا (میتوانیم نشان دهیم پاسخ از کدام سند آمده)	پایین (نمیدانیم مدل از کجا پاسخ را آورده)	تفسیر پذیری
بسیار پایین (چون مدل بر اساس اسناد واقعی پاسخ میدهد)	نسبتاً بالا (مخصوصاً برای اطلاعات خارج از دادههای آموزشی)	توهم
بله (با به روزرسانی پایگاه داده، دانش مدل بهروز میشود)	خیر (دانش مدل ثابت میماند)	دسترسی به دانش جدید
بسیار بالا (میتوانیم ترابایت سند داشته باشیم)	محدود (با افزایش داده، هزینه بالا میرود)	مقیاس پذیری



RAG

Advanced RAG Techniques



Hybrid Retrieval

Combines semantic (vector) search with keyword (BM25) search.

Best for: 'Find exactly Eurocode 7, Table A.3' — exact term matches matter.



Reranking

Retrieve top-20 chunks → cross-encoder re-scores → use top-5.

Result: Much higher precision on complex engineering queries.



Metadata Filtering

Filter before semantic search by document type, project ID, or date range.

Example: 'Only search inspection reports 2020-2023 for Project ID XYZ'.



HyDE (Hypothetical Document)

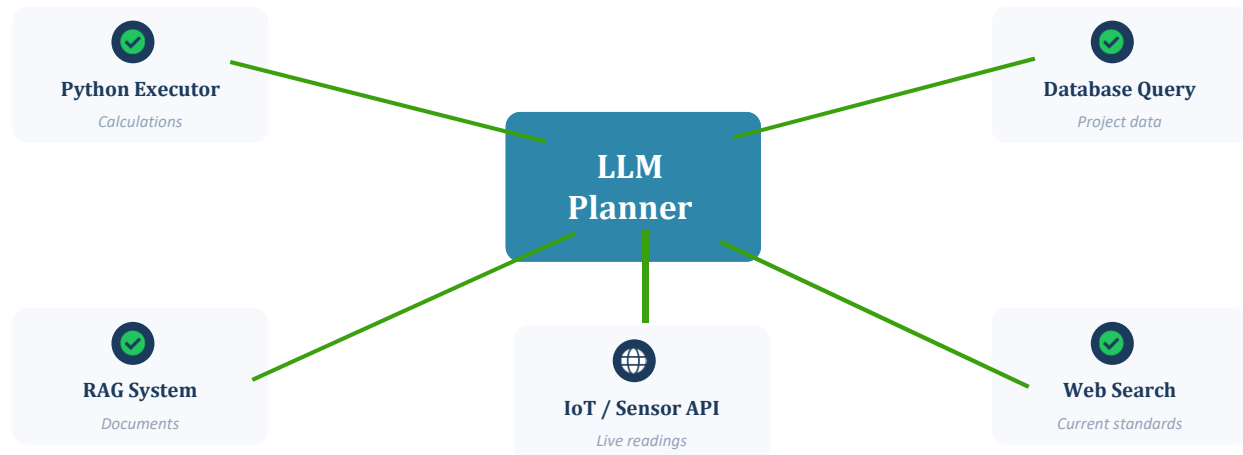
LLM generates a 'hypothetical ideal answer', then retrieve documents similar to it.

Better recall for vague or complex technical queries.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

AGENTS

Agentic LLMs: The Model as Orchestrator

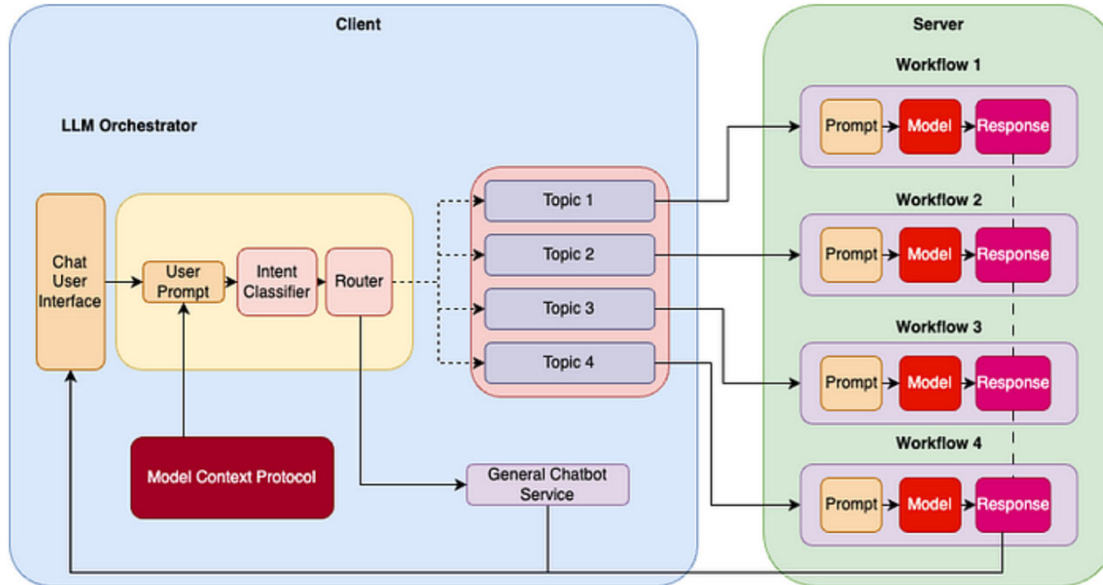


مخزنی از اطلاعات مربوط به هر موضوع خاص درست می‌شود. یک هماهنگ کننده (عامل) قصد کاربر از سوالاتش مشخص می‌شود و هماهنگ کننده کاربر را به مخزن اطلاعات مربوطه وصل می‌کند.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

AGENTS

Agentic LLMs: The Model as Orchestrator



Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

APPLICATION 1

Automated Specification Compliance Checking

The Problem

50 specification clauses × 30-page submittal = hours of manual checking, subject to human error and fatigue.

The AI-Assisted Solution

- 1 Project spec ingested into RAG system
- 2 Submittal parameters extracted by parser
- 3 LLM checks each clause → compliance flag
- 4 Engineer reviews flagged items only

60-70% reduction in review time · Engineers focus only on flagged items

⚠ Output must always be reviewed by a licensed engineer

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

APPLICATION 2

Structural Inspection Report Generation



Inspector Input

Structured field notes · Defect photos · Condition observations



LLM Generates

Standard report format · NBI condition ratings · Critical flagging · Maintenance recs



Engineer Reviews

Validates AI output · Signs off report · Focuses on judgment

Northwestern Univ. 2023: 87% agreement with experienced inspectors

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

APPLICATION 3

Geotechnical Data Intelligence

RAG-powered query example:

Q: "Are there zones where $N < 5$ within 3m of foundation level?"

A: "Yes — BH-03 at 2.5m shows $N=3$ in sandy silt (ML), and BH-11 at 1.8m shows $N=3$ in loose fine sand (SP). Both zones require attention for liquefaction potential at design PGA."

AI Capabilities — Geotechnical Data

- Identify anomalous zones across the site
- Correlate SPT with lab shear strength for QA checks
- Draft Geotechnical Conditions report section

Fugro — world's largest geotech firm — investing in AI-assisted subsurface interpretation

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

APPLICATION 4

Engineering Code Q&A Assistant (RAG-powered)

WITH RAG — Grounded Answer

Q: "Eurocode 2 — min. cover, XC3 exposure, 100-year design life?"

A: "Per EN 1992-1-1 §4.4.1.2, Table 4.4N:

$c_{min,dur} = 25$ mm for XC3, structural class S4.
With 100-year design life modifier ($\Delta c_{DL} = 0$ as Table NA.4), and $\Delta c_{dur,\gamma} = 0$:

$c_{nom} = c_{min} + \Delta c_{dev} = 25 + 10 = 35$ mm

[Source: EN 1992-1-1:2004, Table 4.4N]"



Cites actual document — no hallucination

WITHOUT RAG — Hallucination Risk

Q: Same question

A: "For XC3 exposure with 100-year design life, Eurocode 2 specifies a nominal cover of 40mm for the structural class adjustment..."
[fabricated details — section references may be wrong, values inaccurate]



Plausible-sounding, potentially wrong

RAG retrieves the actual standard text — eliminating the main source of hallucination on code questions.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT |

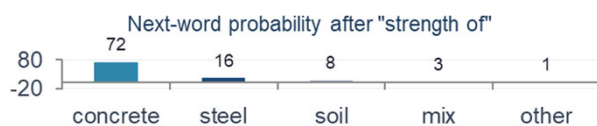
APPLICATION 5

Construction Delay Analysis

The Problem: Thousands of emails, RFIs, meeting minutes & schedule updates to reconstruct the project timeline and assign responsibility.

Multi-Agent System

1. Extract dated events from all correspondence
2. Map cause-effect relationships between delays
3. Attribute responsibility (RFI turnaround times)
4. Generate preliminary delay narrative draft



Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT |

FRAMEWORK

Deployment Decision Framework



Needs domain style?

→ **Fine-Tune**

Needs your documents?

→ **RAG**

Needs calculations?

→ **Tool Use**

Multi-step workflow?

→ **Multi-Agent**

Simple, well-defined Q?

→ **Prompt Engineering**

1. Define task precisely · 2. Build human-in-the-loop review · 3. Validate on known-answer cases · 4. Document the system · 5. Never skip engineer review on safety-critical outputs

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT

ETHICS & LIMITS

What AI Cannot Do — And What You Must Do

Technical Limitations



Hallucination not eliminated — RAG & fine-tuning reduce, not eliminate



Multi-step arithmetic errors — use tool-based code execution



Context window limits — large projects exceed capacity



Outdated knowledge — code updates may not be reflected



Imperfect retrieval — relevant chunks may not be found

AI amplifies engineers — it does not replace engineering judgment.

Prof. Hasan Ghasemzadeh | K.N. Toosi University of Technology | Adv. App. of AI & DT | Session 2

Again The Profession Adapts

"Slide rule → Calculator → FEM → LLMs.
Each time, the technology amplified the engineer — not replaced them."



Understand the tool → earn the right to use it



AI strengthen expert judgment, not replace it



The question is not if — it's how responsibly

پروژه برنامه نویسی

تمرین: یک مساله عمرانی را به روش هوش مصنوعی حل نمایید. برخی عناوین به شرح ذیل هستند

- ✓ پیش بینی ظرفیت باربری شمع با استفاده از داده های CPT و XGBoost
- ✓ خوشه بندی لایه های خاک با استفاده از K-Means روی داده های Cone Penetration Test (CPT)
- ✓ پیش بینی سری زمانی نشست پی سطحی با LSTM (داده های ابزار دقیق)
- ✓ شناسایی الگوی ترک در تصاویر بتن با CNN
- ✓ بهینه سازی طراحی گود عمیق (top-down) با استفاده از الگوریتم ژنتیک (GA)
- ✓ شبیه سازی جریان دانه ای (ریزش ماسه) با GNN پیاده سازی ساده شده مقاله Choi & Kumar 2024
- ✓ ایجاد یک دوقلوی دیجیتال ساده برای یک شیروانی خاکی (نشست لحظه ای VS هشدار)
- ✓ استفاده از مدل زبانی بزرگ برای کنترل دفترچه محاسبات سازه یا ژئوتکنیک



به دنبال آرزوهایم خواهم رفت
مسکرم نمیرم عهد بسته ام قبل از