



رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن باز هوش مصنوعی:
راهنمای اتصال به مدل ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



راهنمای اتصال به مدل ها از طریق API و استفاده از خدمت LLM Serving

بهمن ماه ۱۴۰۴

تاریخ تهیه: ۱۴۰۴/۰۸/۲۵

تاریخ به روزرسانی: ۱۴۰۴/۱۱/۰۵



رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن باز هوش مصنوعی:
راهنمای اتصال به مدل ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



فهرست مطالب

۱- مقدمه	۳
۱-۱ لیست مدل هایی که فعلاً در سکو موجود هستن (به روز شده در تاریخ ۱۴۰۴/۱۱/۰۵)	۳
۲- پیش نیازها	۳
۲-۱ واسط API برای گرفتن لیست مدل های در دسترس	۴
۲-۲ واسط API برای ارسال یک درخواست و دریافت پاسخ	۶
۲-۳ پارامترهای رایج در Chat Completions	۸
۲-۴ تست مدل های Embedding	۹
۲-۵ تست مدل های Reranker	۱۰
۲-۶ تست مدل های Audio (ASR)	۱۳
۲-۷ سایر دستورات API موجود	۱۴
۲-۸ پیاده سازی RAG و استفاده از طریق API	۱۷



راست جمهوری
معاونت علمی و فناوری

سکوی ملی متن‌باز هوش مصنوعی: راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



۱ مقدمه

این راهنما توضیح می‌دهد چگونه می‌توانید با استفاده از API به مدل‌های زبانی (LLM) متصل شوید، درخواست بفرستید و پاسخ دریافت کنید.

۱-۱ لیست مدل‌هایی که فعلاً در سکو موجود هستن (به‌روزرشته در تاریخ ۱۴۰۴/۱۱/۰۵)

تاریخ به‌روزرسانی: ۱۴۰۴/۱۱/۰۵		
اندازه تقریبی پارامتر	نام مدل	نوع مدل
8B	Qwen3-VL-8B-Instruct	LLM
235B	Qwen3-VL-235B-Instruct	
685B	DeepSeek-V3.1	
32B	Qwen3-32B	
120B	gpt-oss-120b	
72B	Qwen2.5-72B	
32B	Qwen3-32B-128K	
17B	Llama4-Scout-17B-16E	
27B	gemma-3-27b-it	
8B	Dorna2-Llama3.1-8B-Instruct	
235B	Qwen3-235B-A22B-Instruct-2507	
7B	VulnLLM-R-7B	
N/A	Persian_Sentence_Embedding_v3	SLM (Embedding)
N/A	Jina -embeddings-v3	
N/A	jina-embeddings-v4	
560M	intfloat-multilingual-e5-large	
8B	Qwen3-Embedding-8B	
N/A	baai-bge-m3	
150~200M	jina-reranker-v2-base-multilingual	SLM (Reranker)
200~300M	baai-bge-reranker-v2-m3	
3B	DeepSeek-OCR	Vision (OCR)
N/A	rednote-hilab-dots.ocr	
1.55B	whisper-large-v3	Audio (ASR)
1.55B	whisper-large-v3-fa	

۲ پیش‌نیازها

۱. توکن دسترسی (API Key)

- مثال: sk-xxxxxxxxxxxxxxxxxxxxxxxx (توکن شخصی خودتان که برای شما ساخته شده است باید جایگزین xها شود)

۲. ابزار ارسال درخواست HTTP

- Postman
- curl
- Python یا JavaScript

۳. آدرس API

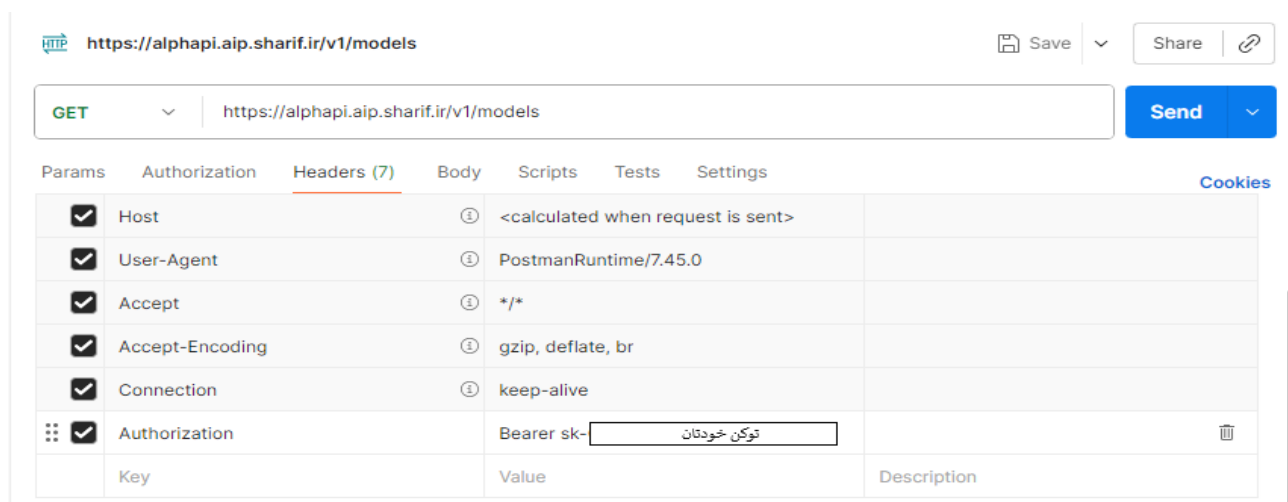
<https://alphapi.aip.sharif.ir>

۲-۱ واسط API برای گرفتن لیست مدل‌های در دسترس

۱. ارسال درخواست از طریق Postman

در Postman برای مطمئن شدن از نام دقیق مدل‌ها باید یک ریکوئست GET بسازید.

- روی New کلیک کنید و یک HTTP Request جدید بسازید.
- متد (Method) را روی GET بگذارید.
- آدرس (URL) را مطابق زیر بنویسید:
- <https://alphapi.aip.sharif.ir/v1/models>
- به تب Headers بروید و توکن خودتان را با عنوان Authorization مشابه تصویر ۱ وارد کنید:
- Authorization → Bearer sk- توکن خودتان
- نیازی به Body نیست.
- روی Send بزنید.



تصویر ۱ تنظیمات احراز هویت در تب Header

خروجی درخواست GET بالا مشابه تصویر ۲ خواهد بود. ID مدل‌هایی که توکن مربوطه به آن‌ها دسترسی دارد قابل رویت است. هر کدام از idها نام دقیق مدلی است که باید در درخواست‌های بعدی استفاده کنید.



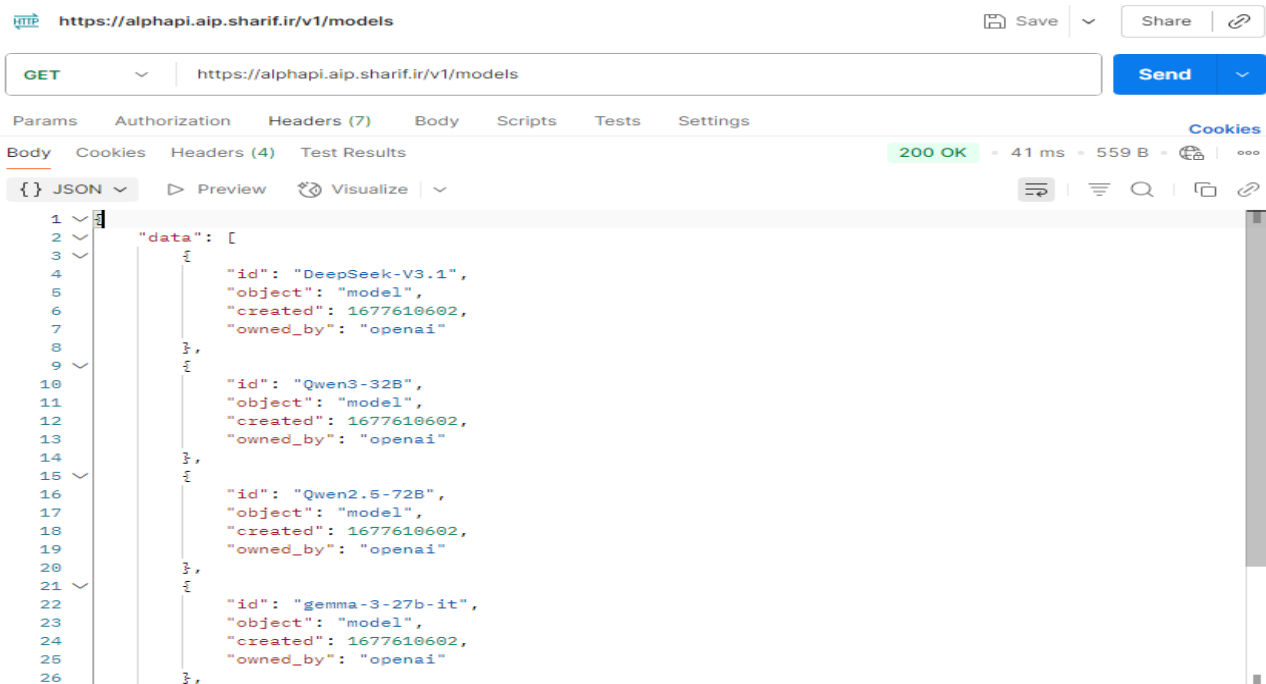
رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن‌باز هوش مصنوعی: راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



تصویر ۲ خروجی درخواست GET برای دیدن نام و مشخصات مدل‌های موجود

۲. ارسال درخواست از طریق Curl

دستور زیر را می‌توان برای دریافت خروجی ارسال کرد:



```
curl.exe -s -X GET "https://alphapi.aip.sharif.ir/v1/models" ^
-H "Authorization: Bearer sk-توکن-شما" |
ConvertFrom-Json | ConvertTo-Json -Depth 10
```

۳. ارسال درخواست از طریق Python

```
import requests
URL = "https://alphapi.aip.sharif.ir/v1/models"
TOKEN = "sk-توکن-شما"
headers = {
    "Authorization": f"Bearer {TOKEN}",
    "Accept": "application/json"
}
resp = requests.get(URL, headers=headers, timeout=30)
resp.raise_for_status()

data = resp.json()

# چاپ فقط نام مدل‌ها
for m in data.get("data", []):
    print(m.get("id"))
```

 ریاست جمهوری معاونت علمی و فناوری	سکوی ملی متن‌باز هوش مصنوعی: راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API			
	پیوست: ندارد	تاریخ: ۱۴۰۴/۱۱/۰۵	شماره: ۰۶	

۲-۲ واسط API برای ارسال یک درخواست و دریافت پاسخ

چهار روش برای اتصال به مدل‌ها و دریافت پاسخ در این مرحله نیز وجود دارد که به ترتیب توضیح داده می‌شود.

۱. ارسال درخواست از طریق Postman

- یک Request جدید بسازید.
- متد URL روی POST باشد و آدرس زیر وارد شود:
- به تب Headers بروید و دو مورد زیر را اضافه کنید:

<https://alphapi.aip.sharif.ir/v1/chat/completions>

Content-Type: application/json

Authorization: Bearer sk - توکن خودتان

- به تب Body بروید، حالت را روی raw و فرمت را روی JSON بگذارید.
- محتوای Body را به صورت امتحانی مشابه تصویر زیر پر کنید.
- روی Send بزنید و پاسخ مدل را در پایین صفحه و در عکس زیر ببینید.



رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن‌باز هوش مصنوعی: راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



<https://alphapi.aip.sharif.ir/v1/chat/completions>

Save

Share

POST

<https://alphapi.aip.sharif.ir/v1/chat/completions>

Send

Params Authorization Headers (9) Body Scripts Tests Settings

Cookies

Beautify

☐ none ☐ form-data ☐ x-www-form-urlencoded ☒ raw ☐ binary ☐ GraphQL JSON

```
1 {
2   "model": "DeepSeek-V3.1",
3   "messages": [
4     {
5       "role": "user",
6       "content": "What is the capital of Iran?"
7     }
8   ]
9 }
10
```



تصویر ۳ نمونه ارسال یک سوال یا Prompt به یکی از مدل‌های سکو

Body Cookies Headers (20) Test Results

200 OK • 2.85 s • 1.71 KB • 🌐 ⋮

{ } JSON

Preview

Visualize

```
1 {
2   "id": "chatcmpl-3f9d7f6c-bb43-4cb4-9c28-c33754925d24",
3   "created": 1757158752,
4   "model": "DeepSeek-V3.1",
5   "object": "chat.completion",
6   "choices": [
7     {
8       "finish_reason": "stop",
9       "index": 0,
10      "message": {
11        "content": "Of course! The capital of Iran is **Tehran**.\n\nHere are some key details about it:\n\n* **Population:** Tehran is Iran's largest city, with a metropolitan population of over 15 million people.\n* **Function:** It is the political, economic, and cultural center of the country, housing the seats of the national government.\n* **Location:** The city is situated in the north-central part of Iran, at the foot of the Alborz mountain range.\n\nTehran became the capital of Iran in 1795 under the Qajar dynasty, Agha Mohammad Khan.",
12        "role": "assistant"
13      },
14      "provider_specific_fields": {
15        "stop_reason": null
16      }
17    }
18  ]
19 }
```

تصویر ۴ خروجی دریافتی از مدل بنا بر سوال ارسالی



رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن‌باز هوش مصنوعی: راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



۲. ارسال درخواست از طریق Curl

```
curl.exe -X POST "https://alphapi.aip.sharif.ir/v1/chat/completions" ^
-H "Content-Type: application/json" ^
-H "Authorization: Bearer sk-توکن-شما" ^
-d "{\"model\": \"DeepSeek-V3.1\", \"messages\": [{\"role\": \"user\", \"content\": \"What is the capital of Iran?\"}]}"
```

۳. ارسال درخواست از طریق کد Python

```
import requests

url = "https://alphapi.aip.sharif.ir/v1/chat/completions"
headers = {
    "Content-Type": "application/json",
    "Authorization": "Bearer sk-توکن-شما"
}
data = {
    "model": "DeepSeek-V3.1",
    "messages": [
        {"role": "user", "content": "What is the capital of Iran?"}
    ]
}
response = requests.post(url, headers=headers, json=data)
print(response.json())
```

۴. ارسال درخواست از طریق JavaScript (Node.js)

```
import fetch from "node-fetch";
const url = "https://alphapi.aip.sharif.ir/v1/chat/completions";
const response = await fetch(url, {
    method: "POST",
    headers: {
        "Content-Type": "application/json",
        "Authorization": "Bearer sk-توکن_خودت"
    },
    body: JSON.stringify({
        model: "DeepSeek-V3.1",
        messages: [
            { role: "user", content: "What is the capital of Iran?" }
        ]
    })
});
const result = await response.json();
console.log(result);
```

درخواست بالا را می‌توانید با سوالات دیگر و پرسش از مدل‌های دیگر نیز تکرار کنید.

۳-۲ پارامترهای رایج در Chat Completions

در کنار فیلدهای اصلی در Body مثل model و messages، می‌توان از پارامترهای زیر برای کنترل نحوه تولید پاسخ توسط مدل استفاده کرد.



راست جمهوری
معاونت علمی و فناوری

سکوی ملی متن‌باز هوش مصنوعی: راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API

شماره: ۰۶ تاریخ: ۱۴۰۴/۱۱/۰۵ پیوست: ندارد



پارامتر	نوع داده	مقدار پیش‌فرض	توضیح
model	string	—	نام دقیق مدلی که از v1/models/ گرفته‌م مثل DeepSeek-V3.1
messages	array	—	لیست پیام‌ها در قالب نقش (role: user, assistant, system) و متن (content)
temperature	float (۰ تا ۲)	1.0	کنترل خلاقیت مدل؛ مقادیر بالاتر = خروجی متنوع‌تر، مقادیر پایین‌تر = خروجی متمرکزتر
top_p	float (۰ تا ۱)	1.0	برای محدود کردن temperature؛ جایگزین روش nucleus sampling تنوع خروجی
max_tokens	integer	بستگی به مدل دارد	حداکثر تعداد توکن خروجی، اگر خیلی کم باشد جواب ناقص می‌آید
n	integer	1	تعداد پاسخ‌هایی که مدل باید همزمان تولید کند
stream	boolean	false	اگر true باشد خروجی به صورت استریم (قطعه‌قطعه) برمی‌گردد
presence_penalty	float (۰ تا ۲)	0	افزایش احتمال تولید موضوعات جدید (جلوگیری از تکرار)
frequency_penalty	float (۰ تا ۲)	0	کاهش احتمال تکرار همان کلمات/عبارات
stop	string یا array	—	تعریف یک یا چند توکن/رشته که با رسیدن به آن تولید متن متوقف شود
user	string	—	شناسه کاربر (برای لاگ یا مدیریت استفاده)، اختیاری است

۲-۴ تست مدل‌های Embedding

۱. ارسال درخواست از طریق Postman

برای دریافت بردار جمله (Sentence Embedding)، مراحل زیر را انجام دهید:

- یک Request جدید در Postman ایجاد کنید.
- متد را روی POST قرار دهید و آدرس زیر را وارد نمایید:

<https://alphapi.aip.sharif.ir/v1/embeddings>

- به تب Headers بروید و دو مورد زیر را اضافه کنید:

Content-Type: application/json

Authorization: Bearer sk - توکن خودتان

- به تب Body بروید، حالت را روی raw و فرمت را روی JSON بگذارید.
- محتوای Body را به صورت امتحانی مشابه زیر پر کنید.

```
{
  "model": "baai-bge-m3",
  "input": [ "تهران پایتخت ایران است" ]
}
```

خروجی دریافتی مشابه تصویر ۵ خواهد بود.



سکوی ملی متن‌باز هوش مصنوعی: راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API



پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶

```

1 {
2   "model": "baai-bge-m3",
3   "input": ["تهران پایتخت ایران است."]
4 }
5

```

```

1 {
2   "model": "baai-bge-m3",
3   "data": [
4     {
5       "embedding": [
6         0.005341342184692621,
7         0.06665994971990855,
8         -0.021812375131249428,
9         0.03208521088491516,
10        -0.02368769235908985,
11        0.0018340141598135233,
12        -0.012208442751944065,
13        -0.010541425466537476,
14        0.0017278861584628062,
15        0.02785841585396366,
16        -0.027514878668448445,
17        0.014965847128498554,
18        -0.03241962566971779,
19        -0.003355756482815686,
20        -0.0048431637316942215,
21        0.000841841974761337,

```

تصویر ۵ خروجی دریافتی از یک مدل Embedding برای تبدیل جمله به بردار

۲. ارسال درخواست از طریق Curl

```

curl.exe -s -X POST "https://alphapi.aip.sharif.ir/v1/embeddings" ^
-H "Content-Type: application/json" ^
-H "Authorization: Bearer sk-توکن-شما" ^
--data-raw "{\"model\": \"baai-bge-m3\", \"input\": [\"تهران پایتخت ایران است.\"]}"

```

۳. ارسال درخواست از طریق کد Python

```

import requests
url = "https://alphapi.aip.sharif.ir/v1/embeddings"
headers = {
    "Content-Type": "application/json",
    "Authorization": "Bearer sk-توکن-شما"
}
data = {
    "model": "baai-bge-m3",
    "input": ["تهران پایتخت ایران است."]
}
print(requests.post(url, headers=headers, json=data).json())



```

۲-۵ تست مدل‌های Reranker

۱. ارسال درخواست از طریق Postman

برای رتبه‌بندی (Rerank) یک فهرست سند بر اساس یک پرسش، مراحل زیر را انجام دهید:

- یک Request جدید در Postman ایجاد کنید.
- متد را روی POST قرار دهید و آدرس زیر را وارد نمایید:

 <p>راست جمهوری معاونت علمی و فناوری</p>	<p>سکوی ملی متن‌باز هوش مصنوعی:</p> <p>راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API</p>			
	پیوست: ندارد	تاریخ: ۱۴۰۴/۱۱/۰۵	شماره: ۰۶	

<https://alphapi.aip.sharif.ir/v1/rerank>

- به تب Headers بروید و دو مورد زیر را اضافه کنید:

Content-Type: application/json

Authorization: Bearer sk- توکن خودتان

- به تب Body بروید، حالت را روی raw و فرمت را روی JSON بگذارید.

- محتوای Body را به صورت امتحانی مشابه زیر پر کنید.

```
{
  "model": "baai-bge-reranker-v2-m3",
  "query": "بهترین رستوران ایتالیایی در تهران کدام است؟",
  "documents": [
    "[ID:D1] رستوران ایتالیایی لاوینا در تهران واقع شده و منوی کامل پیتزا و پاستا دارد",
    "[ID:D2] کافه کتاب تهران در خیابان انقلاب قرار دارد و محیطی آرام برای مطالعه فراهم می‌کند",
    "[ID:D3] رستوران توسکانی در مرکز تهران پیتزا و پاستای باکیفیتی سرو می‌کند",
    "[ID:D4] Many say the best Italian restaurant in Tehran is a cozy spot with authentic pasta and Neapolitan pizza.",
    "[ID:D5] رستوران ایتالیایی ناپولی در اصفهان محبوب است و سرآشپز ایتالیایی دارد",
    "[ID:D6] یک رستوران یونانی در تهران با مزه‌های مدیترانه‌ای فعالیت می‌کند",
    "[ID:D7] جشنواره فیلم فجر هر سال در تهران برگزار می‌شود و ربطی به غذا ندارد",
    "[ID:D8] رستوران ایتالیایی لاوینا در تهران با سرآشپز مهمان آخر هفته‌ها منوی ویژه ایتالیایی ارائه می‌دهد",
    "[ID:D9] یک کافی‌شاپ ایتالیایی در تهران اسپرسو و تیرامیسو سرو می‌کند",
    "[ID:D10] متن طولانی درباره وضعیت آب‌وهوا و مسابقات فوتبال اروپا که هیچ ارتباطی با رستوران ندارد"
  ],
  "top_n": 3,
  "return_documents": true
}
```

خروجی دریافتی مشابه تصویر ۶ خواهد بود.

https://alphapi.aip.sharif.ir/v1/rerank

POST https://alphapi.aip.sharif.ir/v1/rerank

Params Authorization Headers (10) Body Scripts Settings

none form-data x-www-form-urlencoded raw binary GraphQL JSON

```

1 {
2   "model": "jina-reranker-v2-base-multilingual",
3   "query": "بهترین رستوران ایتالیایی در تهران کدام است؟",
4   "documents": [
5     "[ID:01] رستوران ایتالیایی لاوینا در تهران واقع شده و منوع کامل پیتزا و پاستا دارد",
6     "[ID:02] کافه کتاب تهران در خیابان انقلاب قرار دارد و محیطی آرام برای مطالعه فراهم میکند",
7     "[ID:03] رستوران توسکانی در مرکز تهران پیتزا و پاستای باکیفیتی سرو میکند",
8     "[ID:04] Many say the best Italian restaurant in Tehran is a cozy spot with authentic pasta and Neapolitan pizza.",
9     "[ID:05] رستوران ایتالیایی ناپولی در اطلمان محبوب است و سرآشپز ایتالیایی دارد",
10    "[ID:06] یک رستوران یونانی در تهران با مزه‌های مدیترانه‌ای فعالیت میکند",
11    "[ID:07] جشنواره فیلم فجر هر سال در تهران برگزار میشود و ربطی به غذا ندارد",
12    "[ID:08] رستوران ایتالیایی لاوینا در تهران با سرآشپز بهمان آخر هفته‌ها منوع ویژه ایتالیایی ارائه میدهد",
13    "[ID:09] یک کافی‌شاپ ایتالیایی در تهران اسپرسو و شیرامیسو سرو میکند",
14    "[ID:010] متن طولانی درباره وضعیت آب‌وهوا و مسابقات فوتبال اروپا که هیچ ارتباطی با رستوران ندارد"
15  ]
16 }

```

Body Cookies Headers (11) Test Results

JSON Preview Visualize

```

1 {
2   "id": "rerank-7b236b5af46a420fb7df7af20d9952578",
3   "results": [
4     {
5       "index": 3,
6       "relevance_score": 0.8264832602310181,
7       "document": {
8         "text": "[ID:04] Many say the best Italian restaurant in Tehran is a cozy spot with authentic pasta and Neapolitan pizza."
9       }
10    },
11    {
12       "index": 0,
13       "relevance_score": 0.8029361367225647,
14       "document": {
15         "text": "[ID:01] رستوران ایتالیایی لاوینا در تهران واقع شده و منوع کامل پیتزا و پاستا دارد"
16       }
17    },
18    {
19       "index": 7,
20       "relevance_score": 0.6365631810771362,

```

تصویر ۶ خروجی یک مدل Reranker

۲. ارسال درخواست از طریق Curl

```

curl.exe -X POST "https://alphapi.aip.sharif.ir/v1/rerank" ^
-H "Content-Type: application/json" ^
-H "Authorization: Bearer sk-توکن-خودتان" ^
--data-raw "{
  \"model\": \"jina-reranker-v2-base-multilingual\",
  \"query\": \"پایتخت ایران کجاست؟\",
  \"documents\": [
    \"تهران بزرگ‌ترین شهر ایران است و مرکز سیاسی کشور\",
    \"اصفهان یکی از شهرهای تاریخی ایران است\",
    \"تبریز مرکز استان آذربایجان شرقی است\"
  ],
  \"top_n\": 2,
  \"return_documents\": true
}"

```

۳. ارسال درخواست از طریق کد Python

```

import requests

url = "https://alphapi.aip.sharif.ir/v1/rerank"
headers = {
  "Content-Type": "application/json",
  "Authorization": "Bearer sk-توکن-خودتان"
}

```

```
data = {
    "model": "jina-reranker-v2-base-multilingual",
    "query": "پایتخت ایران کجاست؟",
    "documents": [
        "تهران بزرگ‌ترین شهر ایران است و مرکز سیاسی کشور",
        "اصفهان یکی از شهرهای تاریخی ایران است",
        "تبریز مرکز استان آذربایجان شرقی است"
    ],
    "top_n": 2,
    "return_documents": True
}

resp = requests.post(url, headers=headers, json=data)
print(resp.json())
```

۲-۶ تست مدل‌های Audio (ASR)

۱. ارسال درخواست از طریق Postman

برای استفاده از مدل‌های تبدیل صوت به متن همچون Whisper باید مطابق روش زیر عمل کرد. توجه شود که برای استفاده از این دو مدل (Whisper اصلی و Whisper فارسی) باید از آدرس API دیگری استفاده کرد.

- یک Request جدید در Postman ایجاد کنید.
- متد را روی POST قرار دهید و آدرس زیر را وارد نمایید:
- به تب Headers بروید و سه مورد زیر را اضافه کنید:

<https://alphapi.aip.sharif.ir/new/llm/audio/transcriptions>

Authorization: Bearer sk - توکن خودتان

Accept: application/json

Content-Type: multipart/form-data

- به تب Body بروید و گزینه‌ی from-data را انتخاب کنید. فیلدها را نیز مطابق جدول زیر پر نمایید.

Value	Type	Key
فایل صوتی خود (audio.mp3)	File	file
whisper-large-v3-fa	Text	model

- فایل موردنظر با یکی از فرمت‌های mp3 / wav / m4a انتخاب و بارگذاری شود.
- گزینه‌ی send زده شود. پاسخ که متن آن محتوای صوتی بارگذاری شده است، در این مرحله قابل رویت است.

۲. ارسال درخواست از طریق Curl

```
curl -X POST "https://alphapi.aip.sharif.ir/new/llm/audio/transcriptions" ^ -H "accept: application/json" ^ -H "Authorization: Bearer YOUR_KEY" ^ -F "file=@C:\Users\Test\Downloads\audio.mp3;type=audio/mpeg" ^ -F "model=whisper-large-v3"
```

۳. ارسال درخواست از طریق Python

```
import requests
url = "https://alphapi.aip.sharif.ir/new/llm/audio/transcriptions"
```



رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن باز هوش مصنوعی:
راهنمای اتصال به مدل ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



```
headers = {
    "Authorization": "Bearer YOUR_KEY",
    "accept": "application/json",
}

files = {
    "file": ("audio.mp3", open(r"C:\Users\Test\Downloads\audio.mp3", "rb"), "audio/mpeg")
}

data = {
    "model": "whisper-large-v3"
}

response = requests.post(
    url,
    headers=headers,
    files=files,
    data=data
)

print("Status code:", response.status_code)
print("Response:", response.text)
```

۷-۲ سایر دستورات API موجود

خطاهای رایج / تفسیر	Body (نمونه ی حداقلی)	Headers	URL	Method	Feature
۴۰۱: نامعتبر. ۴۲۲/۴۰۰: ساختار Body نام مدل اشتباه.	{ "model": "DeepSeek-V3.1", "messages": [{ "role": "user", "content": "سلام! حالت چطوره؟" }] }	Authorization , Content- Type	https://alphapi.aip.shar if.ir/v1/chat/completi ons	POST	Chat Completion (چت عادی) تولید پاسخ مکالمه ای
۴۰۴/۴۰۰: مدل روی سرور وجود ندارد.	{ "model": "deepseek-r1", "messages": [{ "role": "user", "content": "چرا آسمان آبی است؟" }] }	Authorization , Content- Type	https://alphapi.aip.shar if.ir/v1/chat/completi ons	POST	Reasoning پاسخ تحلیلی/مرحله ای
۴۰۰: فرمت messages اشتباه است	{ "model": "DeepSeek-V3.1", "messages": [{ "role": "system", "content": "تو یک دستیار مفید" هستی." }, { "role": "user", 	Authorization , Content- Type	https://alphapi.aip.shar if.ir/v1/chat/completi ons	POST	Multi-round Conversation حفظ سابقه گفتگو



رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن‌باز هوش مصنوعی:
راهنمای اتصال به مدل‌ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



خطاهای رایج / تفسیر	Body (نمونه‌ی حداقلی)	Headers	URL	Method	Feature
	<pre> "content": " من می‌خوام به تهران سفر "کنم. چه جاهایی رو پیشنهاد می‌دی؟ , { "role": "assistant", "content": " میتونی از بام‌لند، دربند و "کاخ گلستان دیدن کنی , { "role": "user", "content": " برای غذا چی پیشنهاد می‌کنی؟ }] }</pre>				
—	<pre> { "model": "DeepSeek-V3.1", "messages": [{ "role": "system", "content": " فقط با یک جمله کوتاه "جواب بده" }, { "role": "user", "content": " در مورد فواید خواب کافی "بگو." }] }</pre>	Authorization , Content- Type	https://alphapi.aip.shar if.ir/v1/chat/completio ns	POST	Chat Prefix Completion (Beta) وادر کردن مدل به شروع/سبک خاص
۴۰۰: اگر مدل/سرور response_forma t را پشتیبانی نکند.	<pre> { "model": "DeepSeek-V3.1", "response_format": { "type": "json_object" }, "messages": [{ "role": "user", "content": " سه میوه و قیمت تقریبی "آن‌ها را بده" }] }</pre>	Authorization , Content- Type	https://alphapi.aip.shar if.ir/v1/chat/completio ns	POST	JSON Output خروجی JSON معتبر
۴۰۰: ساختار / tools اسامی اشتباه است.	<pre> { "model": "DeepSeek-V3.1", "messages": [{ "role": "user", "content": " هوا در تهران چگونه؟" }], "tools": [{ "type": "function", "function": { "name": "get_weather", </pre>	Authorization , Content- Type	https://alphapi.aip.shar if.ir/v1/chat/completio ns	POST	Function Calling انتخاب/فراخوانی تابع توسط مدل



رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن باز هوش مصنوعی:
راهنمای اتصال به مدل ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



خطاهای رایج / تفسیر	Body (نمونه ی حداقلی)	Headers	URL	Method	Feature
	<pre>"description": " دریافت وضعیت هوا ", "parameters": { "type": "object", "properties": { "city": { "type": "string" } }, "required": ["city"] } }, "tool_choice": "auto" }</pre>				
۴۰۵/۴۰۴: اگر این مسیر غیرفعال باشد.	<pre>{ "model": "DeepSeek-V3.1", "prompt": "Translate to French: Hello world", "max_tokens": 64 }</pre>	Authorization , Content-Type	https://alphapi.aip.sharif.ir/v1/completions	POST	Completions (متن ساده) تکمیل متن غیر چتی
۴۰۰: اگر prompt/suffix ناقص باشد.	<pre>{ "model": "DeepSeek-V3.1", "prompt": "function add(a, b) {", "suffix": "}", "max_tokens": 64 }</pre>	Authorization , Content-Type	https://alphapi.aip.sharif.ir/v1/completions	POST	FIM Completion (Fill-in-the-Middle) پر کردن وسط متن/اکد
—	<pre>{ "model": "DeepSeek-V3.1", "messages": [{ "role": "user", "content": "این متن نسبتاً بزرگ است. bullet خلاصه کن." }], "temperature": 0.2, "context_caching": true }</pre>	Authorization , Content-Type	https://alphapi.aip.sharif.ir/v1/chat/completions	POST (دو بار پیاپی)	Context Caching کش سمت سرور برای ورودی های تکراری
۴۰۱: توکن نامعتبر.	بدون Body	Authorization (+اختیاری: Accept: application/json)	https://alphapi.aip.sharif.ir/v1/models	GET	List Models دیدن مدل های در دسترس
۴۰۴/۴۰۱: مسیر غیرفعال/مجوز ندارید.	بدون Body	Authorization	https://alphapi.aip.sharif.ir/v1/user/balance	GET	Get User Balance دیدن موجودی
400 با پیام هایی مثل provider/model not provided: یعنی مدل embedding روی	<pre>{ "model": "baai-bge-m3", "input": ["This is a test"] }</pre>	Authorization , Content-Type	https://alphapi.aip.sharif.ir/v1/embeddings	POST	Embeddings بردارسازی متن برای مدل هایی که Embedding را ساپورت می کنند.



رأست جمهوری
معاونت علمی و فناوری

سکوی ملی متن باز هوش مصنوعی: راهنمای اتصال به مدل ها و دریافت پاسخ از طریق API

پیوست: ندارد

تاریخ: ۱۴۰۴/۱۱/۰۵

شماره: ۰۶



```
const CTX=scored.map(d=>` - [${d.id}||"doc"} | ${d.source}||"src"}] ${d.snippet}`).join("\n");
pm.environment.set("CTX",CTX);

const body={
  model:MODEL_ID,
  temperature:0.2,
  messages:[
    {role:"system",content:"فقط بر اساس متن های زیر پاسخ بده. اگر کافی نبود بگو اطلاعات کافی نیست"},
    {role:"user",content:`سؤال: ${QUESTION}\n\nزمینه: \n${CTX}\n\nپاسخ: `}
  ]
};
pm.request.headers.upsert({key:"Content-Type",value:"application/json"});
pm.request.body.update(JSON.stringify(body,null,2));
```

۶. محتوای Body را مشابه زیر پر نمایید:

```
{
  "id": "HR-Leave",
  "source": "hr_wiki.md",
  "text": "ثبت شود. تایید مدیر مستقیم لازم است HR درخواست مرخصی باید در پورتال"
},
{
  "id": "Procurement",
  "source": "procurement.pdf",
  "text": "انجام می شود PR خرید تجهیزات فقط از طریق واحد تدارکات و ثبت فرم"
},
{
  "id": "Onboarding",
  "source": "onboarding.docx",
  "text": "فرایند استخدام شامل مصاحبه فنی و منابع انسانی است"
},
{
  "id": "Leave-Types",
  "source": "hr_policy.md",
  "text": "مرخصی استعلاجی با ارائه گواهی پزشک تایید می شود. مرخصی ساعتی با تایید مدیر ممکن است"
}
]
```

۱۱. هوشمند با Embedding از طریق Postman

برای استفاده از RAG به صورت کاملاً هوشمند یا افزودن دیتابیس به آن، نیاز است که از مدل های Embedding موجود در سکوی استفاده شود تا ایندکس های برداری ساخته شوند یا دیتابیس برداری شده (Vector DB) به مدل داده شود و سوال هم توسط یک سرویس محلی برداری شود و نتیجه بر اساس آن حاصل شود.