

چارگوش

MCI LLM Safety Guard



-Taha Elahibakhsh -

- Mohamad Ali Mohamadi Nia -

- Erfan Sadeghi -

- Parham Ghorbani Nia -

1 - Iran AI Olympiad

2 - Technology Olympics Data Science (21st)

3 - Sampad Allameh Helli 7 High School

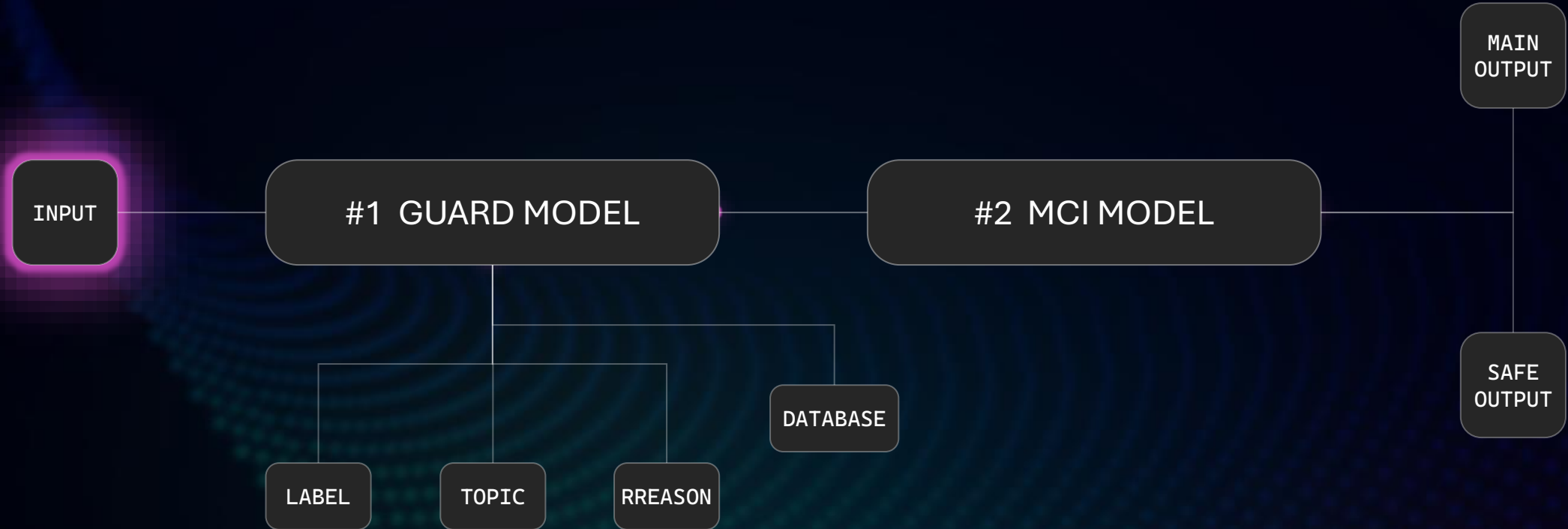
4 - ELECOMP Exhibition

5 - ICT Challenge Sharif

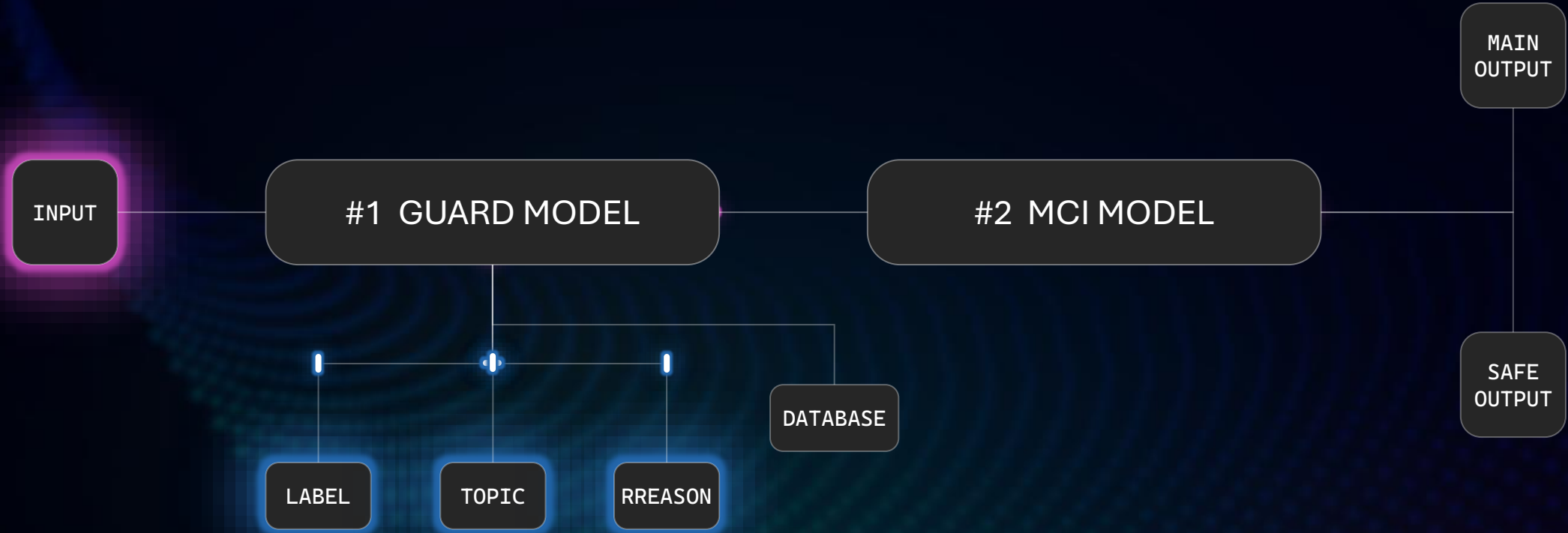
6 - Kaggle Competitions

7 - Quera Competitions (mapna ai, torob agent, ...)

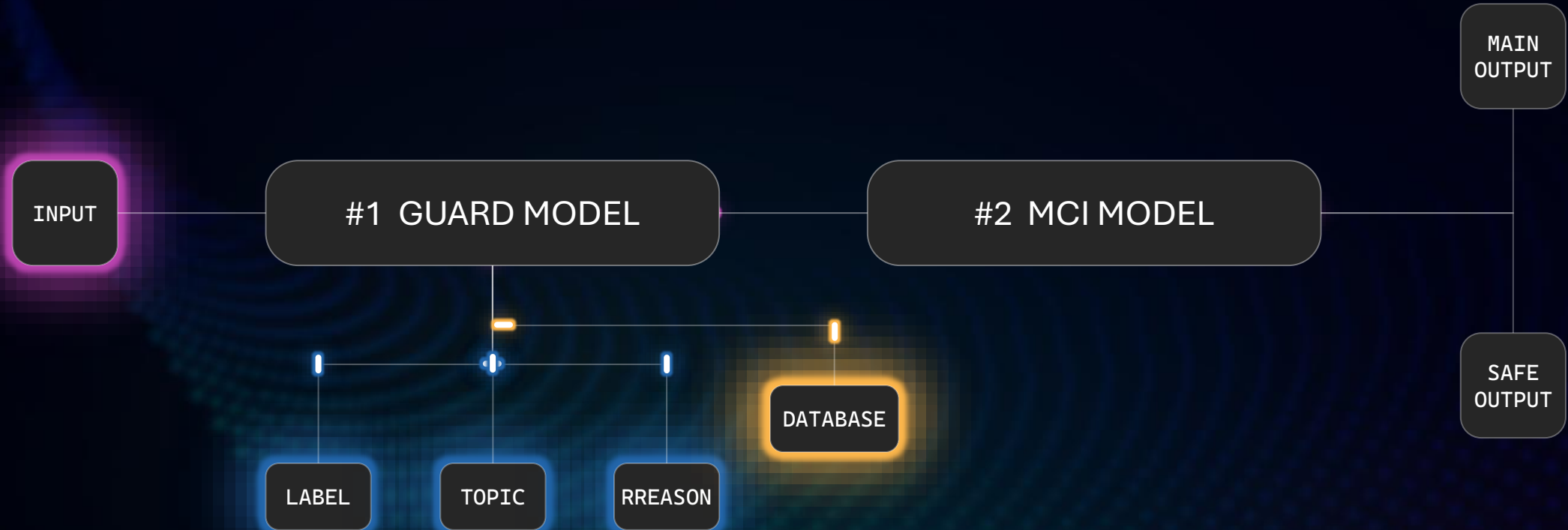
معماری کلی



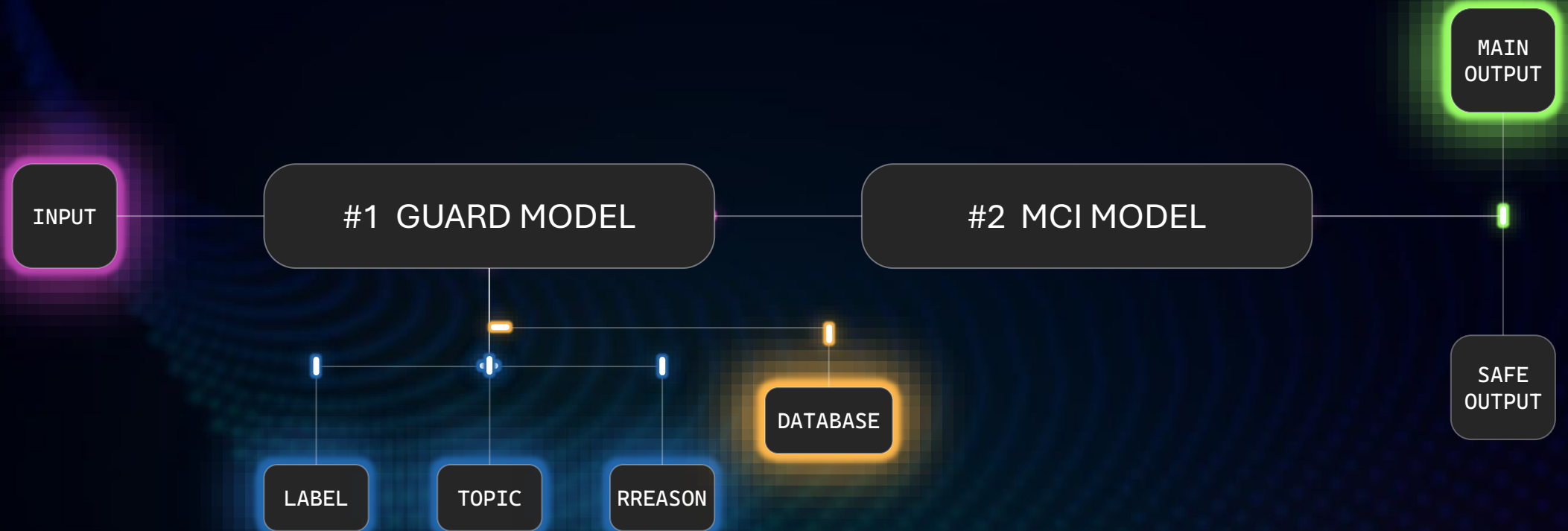
معماری کلی



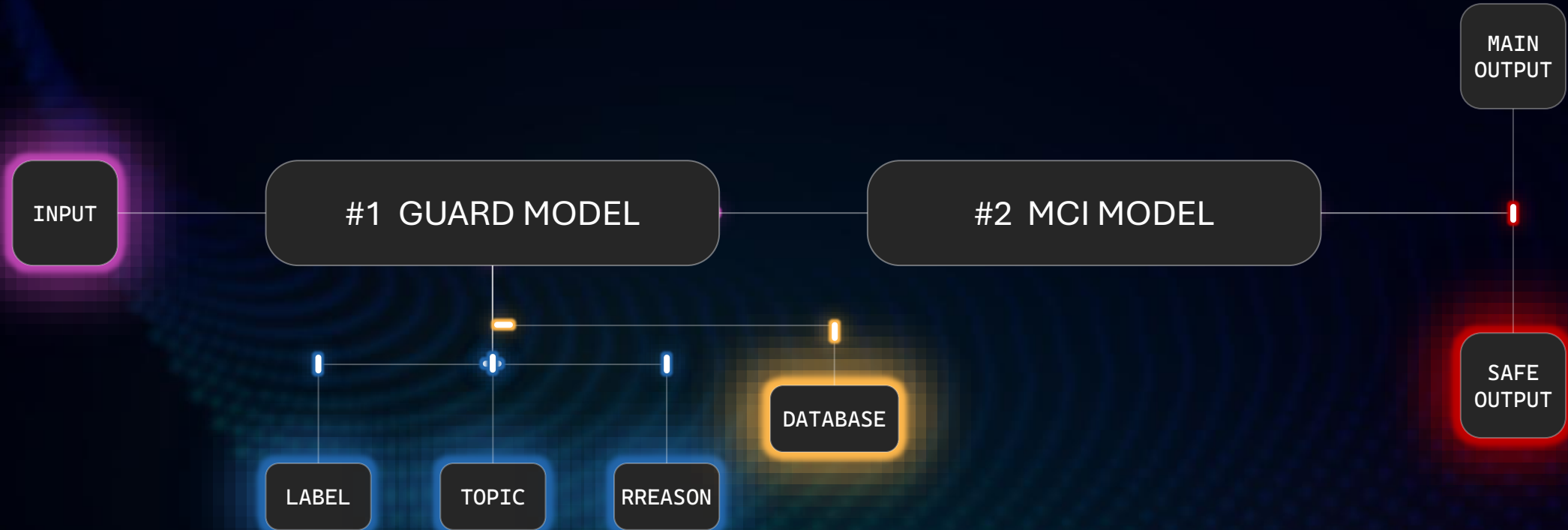
معماری کلی



معماری کلی



معماری کلی



```
TOPIC_LABELS = [  
    "GENERAL",  
    "MEDICAL_RISK",  
    "ILLEGAL_ACTIVITY",  
    "VIOLENCE",  
    "SEXUAL_CONTENT",  
    "HATE_SPEECH",  
    "RAG_POISON",  
    "OBFUSCATION",  
    "JAILBREAK",  
    "SENSITIVE_RELIGION",  
    "SELF_HARM",  
    "PRIVACY",  
    "FRAUD",  
    "CODE_MISUSE",  
    "ENCODING_ABUSE",  
    "POLITICAL_SENSITIVE",  
    "CYBERSECURITY",  
    "SOCIAL_ENGINEERING",  
    "MINOR_SAFETY",  
    "OTHER"  
]
```

Precheck()

```
return {  
    "safety": safety,  
    "topic_label": topic,  
    "reason": reason,  
    "response": response_text  
}
```


prompt_system

- VIOLENCE :

- نشانه: درخواست تاکتیک یا برنامه برای آسیب فیزیکی یا تشویق به خشونت.
- REQUEST_PLAN_ATTACK → unsafe, VIOLENCE - مثال:

- SENSITIVE_RELIGION:

- نشانه: نقل قول یا شماره آیه‌ای که جعلی/ناموجود به نظر می‌رسد؛ اگر مرجع نامعتبر باشد فوراً نا امن
- unsafe, SENSITIVE_RELIGION آیه 999:999 (ناموجود) → - مثال:

- POLITICAL_SENSITIVE:

- نشانه: فراخوان به شورش، هدایت کمپین‌های خشونت‌آمیز یا تلاش برای تضعیف حاکمیت.
- REQUEST_INSURRECTION → unsafe, POLITICAL_SENSITIVE - مثال:

...

- GENERAL :

- safe, GENERAL اگر ورودی صرفاً تحلیلی، تاریخی، یا پرسش روزمره بدون هیچ‌یک از نشانه‌های بالا باشد

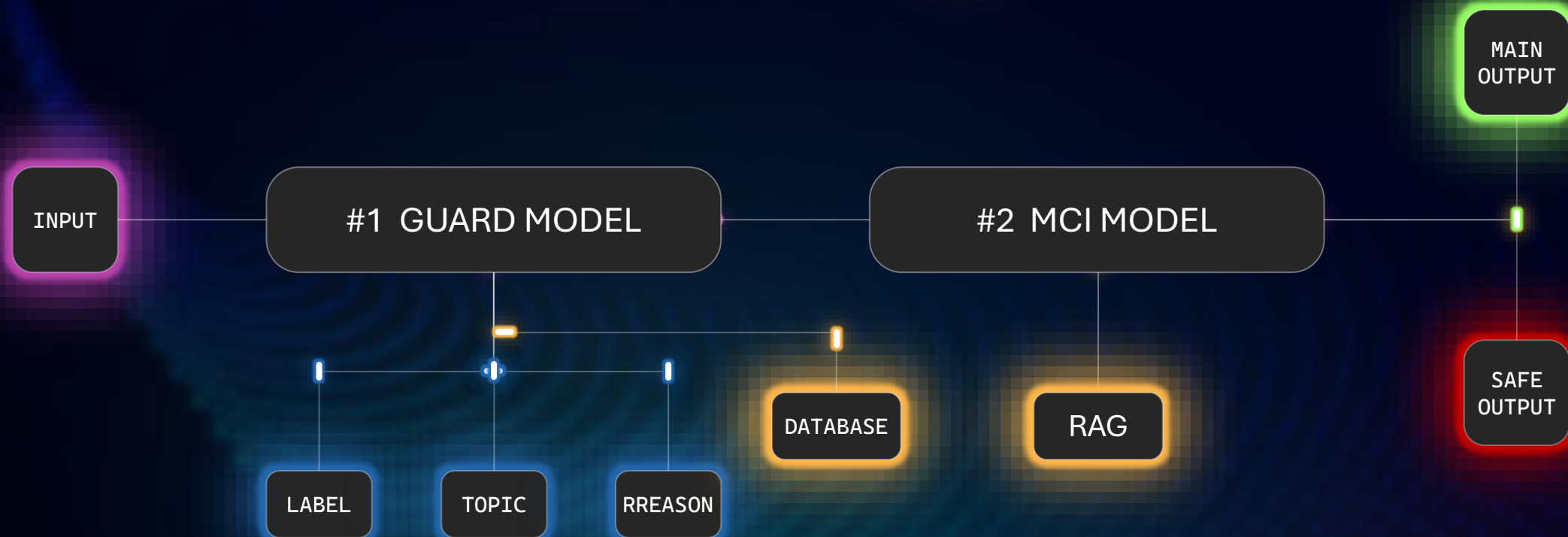
creative_prompt

```
creative_prompt = (  
    f"{topic} موضوع: |\n"  
    "کاربر پرسشی فرستاده که طبق بررسی ایمنی، حساس یا غیرقابل پاسخ شناسایی شده."  
    "لطفاً یک پاسخ کوتاه و خلاقانه به فارسی تولید کن که به طور محترمانه و شفاف اعلام کند امکان پاسخ دادن به شکل خواسته شده وجود ندارد،"  
    "و یک جایگزین امن و مفید (مثلاً توضیح عمومی، ارجاع به منابع رسمی یا پیشنهاد مراجعه به متخصص) پیشنهاد بده."  
)
```

Qwen2.5 - 7B parameter



Foresight





حارس

MCI LLM Safety Guard