

به نام خدا



دانشگاه شهید بهشتی

گزارش پیشرفت فاز سوم پروژه معماری نرم افزار: معماری سیستم های هوش مصنوعی

استاد راهنما درس: دکتر صادق علی اکبری

نام اعضا(تک نفره) : حسن رکن آبادی - شماره دانشجویی : ۴۰۳۴۴۳۲۰۴

خرداد ۱۴۰۴

فهرست مطالب

بخش ۱: تحلیل معماری سیستم های استنتاج هوش مصنوعی (مرور ادبیات) (رویکرد پژوهشی)

- ۱.۱ مقدمه: از مدل تا سرویس؛ چالش های معماری در استقرار هوش مصنوعی
- ۱.۲ الگوهای معماری و تصمیمات کلیدی در انتخاب سکوی استقرار
- ۱.۳ اصول طراحی برای بهینه سازی در سطح مدل و سیستم
- ۱.۴ جمع بندی: درس های معماری برای استقرار سیستم های هوش مصنوعی

بخش ۲: ارزیابی و تحلیل ابزارهای صنعتی در معماری استقرار (رویکرد صنعتی)

- ۲.۱ مقدمه: از تئوری تا عمل؛ انتخاب ابزار به مثابه یک تصمیم معماری
- ۲.۲ ایزوله سازی و تکرارپذیری با Docker
- ۲.۳ زبان مشترک اکوسیستم: استانداردسازی با ONNX
- ۲.۴ موتور بهینه سازی: دستیابی به اوج کارایی با NVIDIA TensorRT
- ۲.۵ هسته مرکزی استقرار: مدیریت و مقیاس پذیری با NVIDIA Triton
- ۲.۶ جمع بندی تحلیل صنعتی

بخش ۳: پیاده سازی و ارزیابی معماری پیشنهادی

- ۳.۱ مقدمه: از طراحی تا واقعیت
- ۳.۲ تحلیل مقایسه ای معماری: دو روش برای استقرار مدل
- ۳.۲.۱ معماری پایه: یک روش ساده اما شکننده
- ۳.۲.۲ معماری پیشنهادی: یک ساختار هوشمند و بهینه
- ۳.۳ تحلیل نتایج عملکردی: مقایسه مستقیم معماری ها
- ۳.۴ جمع بندی فاز پیاده سازی

بخش ۴: معرفی ابزار مقایسه گر و تحلیل بصری نتایج (پیشنهادی)

- ۴.۱ معرفی اپلیکیشن "AI Deployment Architecture Comparator"
- ۴.۲ مقایسه سریع سرویس (Quick Comparison)
- ۴.۳ نتایج تست بار (Load Benchmark)

بخش ۵: جمع بندی نهایی

بخش ۶: منابع مورد استفاده در پروژه

بخش ۱: تحلیل معماری سیستم های استنتاج هوش مصنوعی (مرور ادبیات)

۱.۱ مقدمه: از مدل تا سرویس؛ چالش های معماری در استقرار هوش مصنوعی

در سال های اخیر، کانون توجه در حوزه هوش مصنوعی از چالش صرف «ساخت مدل» به چالش پیچیده تر «عملیاتی سازی و استقرار مدل» تغییر یافته است. امروزه، ساختن یک مدل با دقت بالا، اگرچه دستاوردی مهم است، اما تنها نقطه شروع یک مسیر مهندسی محسوب می شود. چالش اصلی، تبدیل این مدل ایستا به یک سرویس نرم افزاری زنده، پایدار و کارآمد است که بتواند در مقیاس بزرگ به کاربران خدمت رسانی کند. این گذار، یک مسئله عمیق در حوزه معماری نرم افزار است، زیرا محصول نهایی باید مجموعه ای از مشخصه های کیفی بنیادین را برآورده سازد.

از مهم ترین این مشخصه ها می توان به کارایی اشاره کرد که خود به دو معیار کلیدی تقسیم می شود: زمان پاسخ دهی که نشان دهنده سرعت سیستم در پاسخ به یک درخواست منفرد است و توان عملیاتی که ظرفیت پردازش سیستم (مثلاً تعداد تصاویر پردازش شده در ثانیه) در یک بازه زمانی مشخص را معین می کند. در کنار کارایی، مقیاس پذیری^۱ قرار دارد که به توانایی معماری در مدیریت بهینه بارهای کاری رو به رشد و افزایش تعداد کاربران همزمان اشاره دارد. همچنین،

بهینگی در مصرف منابع، به ویژه در استفاده از پردازنده های گرافیکی (GPU) گران قیمت، نقشی حیاتی در کنترل هزینه های عملیاتی ایفا می کند؛ به طوری که یک معماری بهینه می تواند هزینه پردازش یک حجم کاری مشخص را تا ده ها برابر کاهش دهد. در نهایت، قابلیت اصلاح و نگهداری تضمین می کند که سیستم بتواند به سادگی با مدل های جدید به روز رسانی شده و در طول زمان تکامل یابد.^۴ این بخش از گزارش، با بررسی عمیق پژوهش های اخیر، به تحلیل الگوهای معماری، تصمیمات کلیدی و مصالحه های (Trade-offs) موجود در طراحی سیستم های استنتاج می پردازد.

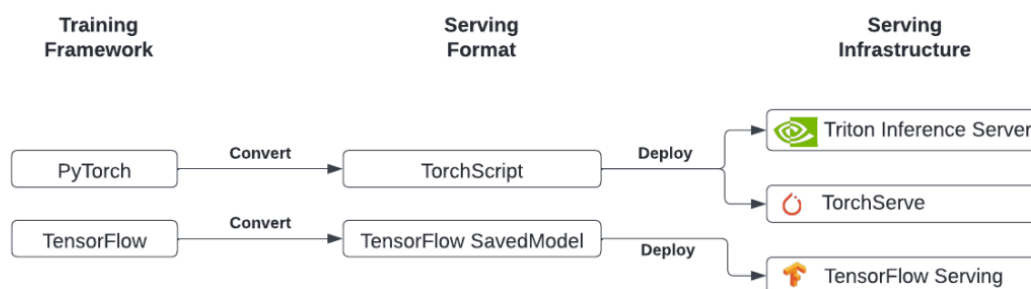
کلمات کلیدی: معماری سیستم های هوش مصنوعی - معماری استقرار - مقیاس پذیری - توان عملیاتی

¹ Scalability

۱.۲ الگوهای معماری و تصمیمات کلیدی در انتخاب سکوی استقرار (رویکرد پژوهشی)

در میان راهکارهای گوناگون، یک الگوی معماری غالب با عنوان «استنتاج به عنوان سرویس»^۲ وجود دارد. این الگو که نمونه‌ای تخصصی از معماری سرویس‌گرا است، بر اصل جداسازی^۳ منطق سنگین و محاسباتی استنتاج از برنامه‌های کاربردی اصلی تأکید دارد. در این ساختار، سرورهایی متمرکز و مجهز به GPU، وظیفه اجرای مدل‌ها را بر عهده گرفته و سایر بخش‌های نرم‌افزار از طریق یک رابط برنامه‌نویسی کاربردی (API) تحت شبکه با آن‌ها ارتباط برقرار می‌کنند. یک پژوهش بسیار روشنگر در این زمینه، کاربرد این الگو را در مراکز محاسبات علمی اشتراکی، که اغلب با کمبود منابع GPU مواجه هستند، مورد مطالعه قرار داده است. نتایج نشان داد که استقرار یک سرور مرکزی مبتنی بر سکوی Nvidia Triton، زمان مورد نیاز برای تحلیل داده‌ها با مدل‌های پیچیده شبکه‌های عصبی گراف (GNN) را تا پنجاه برابر کاهش داده است.^۴ این خود گواهی بر قدرت این الگو در متمرکزسازی و بهینه‌سازی منابع است.

همچنین، برای حل چالش سازگاری محیط‌ها و مدیریت وابستگی‌ها، مقالات پژوهشی به طور گسترده بر نقش محوری فناوری کانتینرسازی با ابزار داکر^۴ تأکید ورزیده‌اند. این فناوری به معمار سیستم اجازه می‌دهد تا هر جزء از سامانه (شامل مدل، سرور استنتاج و کتابخانه‌های مورد نیاز) را در یک محیط ایزوله و قابل تکرار به نام کانتینر بسته‌بندی کند که برای انجام آزمون‌های مقایسه‌ای معتبر، امری ضروری است.^۵



(دیاگرام ۱: معماری مقایسه‌ای سکوی استقرار)

توضیح دیاگرام ۱: این دیاگرام باید به صورت بصری، مسیر یک مدل را از چارچوب آموزشی اولیه آن (مانند PyTorch یا TensorFlow) تا استقرار نهایی روی سه سکوی اصلی (Triton Inference Server، TensorFlow Serving، و TorchServe) نمایش می‌دهد. هدف این تصویر، ارائه یک نمای کلی و مقایسه‌ای از اکوسیستم‌ها و جریان‌های کاری مختلف است تا به درک بهتر تفاوت‌های ساختاری آن‌ها کمک کند.

انتخاب سکوی استقرار، یک تصمیم معماری بنیادین با بده‌بستان‌های مشخص است:

- Nvidia Triton Inference Server: با معماری چند-فریم‌ورکی خود که برای دستیابی به حداکثر کارایی و توان عملیاتی طراحی شده، قدرتمندترین گزینه در محیط‌های مبتنی بر GPU شناخته می‌شود. اما این کارایی به قیمت پیچیدگی بیشتر در راه‌اندازی (نیاز به درایورهای خاص انویدیا و Nvidia Container Toolkit) و نیاز به یک دوره «گرم شدن»^۵ طولانی‌تر برای رسیدن به اوج عملکرد همراه است.

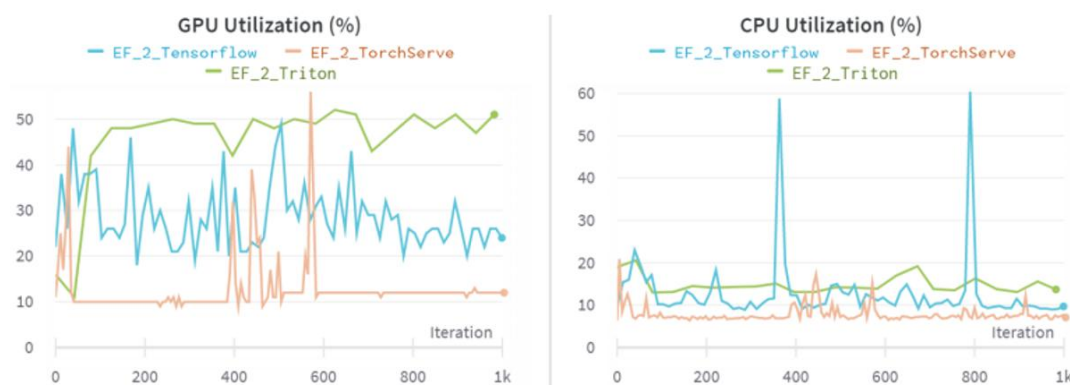
^۲ Inference-as-a-Service

^۳ Decoupling

^۴ Docker

^۵ Warmup

- TensorFlow Serving (TF Serving): برای سادگی و یکپارچگی با اکوسیستم TensorFlow طراحی شده و راه اندازی آن ساده تر است. اما یک بدهستان معماری در آن وجود دارد: قابلیت های پیش پردازش و پس پردازش به صورت داخلی تعبیه نشده اند و منطق آن باید در سمت کلاینت پیاده سازی شود که این امر معماری کلاینت را پیچیده تر می کند.
- TorchServe: این سکو که برای اکوسیستم PyTorch طراحی شده، با ارائه قابلیت بسته بندی کدهای پیش پردازش در قالب یک «پردازشگر» (Handler) در کنار مدل، معماری سمت کلاینت را به مراتب ساده تر می کند. با این حال، این سادگی به قیمت فدا کردن بخشی از کارایی تمام می شود و آزمون های مقایسه ای نشان داده اند که در توان عملیاتی نمی تواند با Triton رقابت کند.



(نمودار ۲: مقایسه مصرف منابع و کارایی)

توضیح نمودار ۲: این نمودار باید نتایج یک آزمون مقایسه ای را در خصوص مصرف منابع (CPU و GPU) و توان عملیاتی سه سکوی مذکور نمایش دهد. این مقایسه بصری به درک بهتر مصالحه میان مصرف منابع و کارایی کمک می کند و نشان می دهد که چگونه سکوی Triton با بهره برداری بسیار بیشتر از GPU (تا ۵۰٪ utilization)، به توان عملیاتی بالاتری نسبت به دو سکوی دیگر (که utilization کمتری دارند) دست می یابد.

۱.۳ اصول طراحی برای بهینه سازی در سطح مدل و سیستم

پس از انتخاب سکو، معمار سیستم می تواند از مجموعه ای از اصول و تاکتیک های طراحی برای بهینه سازی بیشتر بهره گیرد. یک الگوی طراحی مشترک که در پژوهش ها تکرار می شود،

«زنجیره بهینه سازی»^۶ است که معمولاً از مسیر Framework → ONNX → TensorRT پیروی می کند.

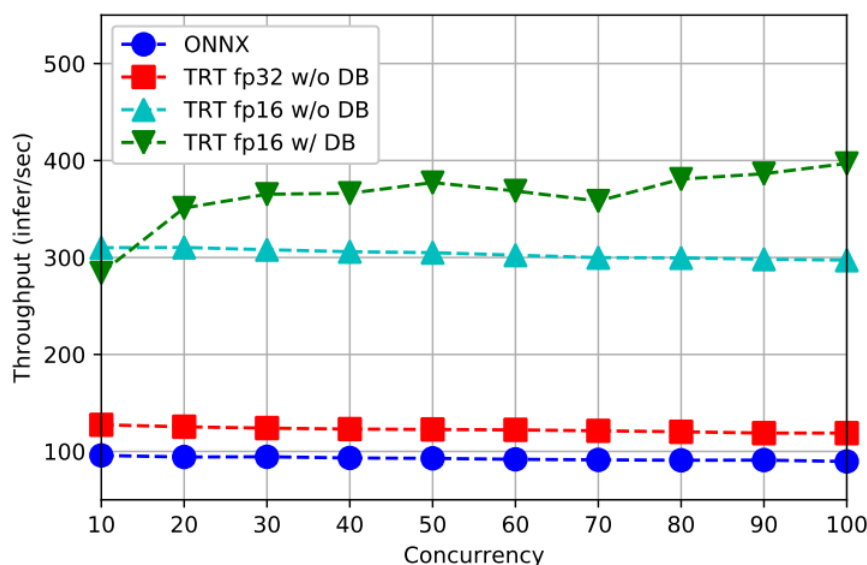
- در این الگو، ONNX به عنوان یک قالب تبادل استاندارد، نقش یک لایه واسطه یا "زبان مشترک" را ایفا می کند که قابلیت همکاری^۷ بین ابزارهای مختلف را ممکن می سازد.
- در مرحله بعد،

Nvidia TensorRT به عنوان یک کامپایلر تخصصی برای سخت افزار انویدیا، مدل را با تکنیک هایی نظیر ادغام لایه ها (Layer Fusion) و انتخاب خودکار بهترین هسته (Kernel Auto-tuning) برای پردازنده گرافیکی مقصد بهینه می کند.

^۶ Optimization Chain

^۷ Interoperability

در کنار این زنجیره، تکنیک کمی سازی یا کوانتیزه سازی^۸ با کاهش دقت عددی محاسبات (برای مثال، از FP32 به FP16)، سرعت پردازش را افزایش داده و حجم مدل را کاهش می دهد. اما یک درس معماری مهم از مقالات این است که این تاکتیک یک راه حل همگانی نیست؛ یک پژوهش نشان داد که کمی سازی در سکوی TorchServe عملاً منجر به کاهش کارایی شده است، در حالی که در Triton و با بهینه سازی TensorRT، بهبود چشمگیری در سرعت ایجاد می کند. در نهایت، در سطح سیستم، تکنیک دسته بندی پویا (Dynamic Batching) که در سکوی Triton پیاده سازی شده، با گروه بندی هوشمند درخواست ها، بهره وری از پردازنده گرافیکی را به حداکثر رسانده و توان عملیاتی را به شدت افزایش می دهد. با این حال، همانطور که در مطالعه استقرار مدل YOLOv5 مشاهده شد، بهره برداری کامل از این قابلیت ممکن است نیازمند توسعه افزونه های سفارشی برای بخش هایی از مدل (مانند اپراتور NMS) باشد. این موضوع نشان دهنده تأثیر متقابل و گاه پیچیده تصمیمات معماری در لایه های مختلف یک سیستم است.



(نمودار ۳: تأثیر انباشتی بهینه سازی ها بر عملکرد)

توضیح نمودار ۳: این نمودار باید به صورت مرحله ای، تأثیر هر تکنیک بهینه سازی را بر مشخصه های کلیدی کارایی مانند زمان پاسخ دهی و توان عملیاتی نشان دهد. به عنوان مثال، از یک خط پایه (مدل در قالب ONNX) شروع کرده و به ترتیب تأثیر افزودن بهینه سازی با TensorRT (حالت FP32)، سپس اعمال کمی سازی FP16 و در نهایت فعال سازی دسته بندی پویا را نمایش می دهد. این تصویر به خوبی ارزش افزوده هر مرحله از زنجیره بهینه سازی را روشن می سازد.

۱.۴ جمع بندی: درس های معماری برای استقرار سیستم های هوش مصنوعی

تحلیل عمیق پژوهش های این حوزه، چندین درس کلیدی و راهبردی را برای یک معمار نرم افزار به ارمغان می آورد. ۱. هیچ معماری واحد و بهینه ای وجود ندارد: انتخاب نهایی همواره نتیجه یک مصالحه هوشمندانه میان مشخصه های کیفی متضاد مانند کارایی، پیچیدگی راه اندازی، هزینه و قابلیت نگهداری است. سکوی Triton برای حداکثر سرعت در محیط های GPU مناسب است، در حالی که TorchServe سادگی کلاینت را در اکوسیستم PyTorch اولویت می دهد.

⁸ Quantization

۲. معماری یک امر چندلایه و یکپارچه است: دستیابی به عملکردی درخشان، نه حاصل یک انتخاب، بلکه نتیجه تعامل و هماهنگی میان لایه سکوی استقرار (مانند Triton)، لایه بهینه‌سازی مدل (مانند TensorRT) و لایه پیکربندی‌های سطح سیستم (مانند Dynamic Batching) است.

۳. اهمیت حیاتی بنچمارک و اعتبارسنجی عملی: تمام مقالات به طور مداوم بر اهمیت حیاتی آزمون‌های مقایسه‌ای و ساخت نمونه‌های اولیه تأکید دارند. مزایای نظری یک تکنیک یا ابزار باید از طریق آزمایش‌های عملی و اندازه‌گیری دقیق در محیط واقعی پروژه اعتبارسنجی شوند، زیرا نتایج ممکن است برخلاف انتظار باشد (مانند تأثیر منفی کوانتیزه‌سازی در TorchServe).

این مبانی و اصول پژوهشی، شالوده علمی و فنی لازم را برای بخش‌های بعدی این پروژه فراهم می‌آورد که در آن‌ها به تحلیل عملی ابزارها و در نهایت، پیاده‌سازی یک معماری بهینه خواهیم پرداخت.

بخش ۲: ارزیابی و تحلیل ابزارهای صنعتی در معماری استقرار (رویکرد صنعتی)

۲.۱ مقدمه: از تئوری تا عمل؛ انتخاب ابزار به مثابه یک تصمیم معماری

پس از آنکه در بخش نخست، الگوهای معماری و مبانی نظری سیستم‌های استنتاج هوش مصنوعی مورد بررسی قرار گرفت، این بخش به تحلیل و ارزیابی عملی ابزارها و فناوری‌های صنعتی می‌پردازد که برای تحقق معماری بهینه در این پروژه انتخاب شده‌اند. در مهندسی نرم‌افزار مدرن، انتخاب یک ابزار صرفاً یک تصمیم فنی نیست، بلکه یک تصمیم معماری بنیادین است که مستقیماً بر ویژگی‌های کیفی سیستم، از جمله کارایی، مقیاس‌پذیری، قابلیت نگهداری و هزینه عملیاتی تأثیر می‌گذارد.

رویکرد این پروژه در انتخاب ابزارها بر پایه یک اصل کلیدی استوار است: استفاده از یک زنجیره ابزار^۹ یکپارچه و اثبات‌شده در صنعت که هر جزء آن، مسئولیت یک بخش مشخص از خط لوله استقرار را بر عهده گرفته و با سایر اجزا تعامل بهینه دارد. این رویکرد در تضاد با انتخاب ابزارهای پراکنده و نامرتب است که گرچه ممکن است هر یک در کار خود بهترین باشند، اما اتصال آن‌ها به یکدیگر هزینه‌های مهندسی و ریسک‌های عملیاتی بالایی را تحمیل می‌کند. ابزارهای منتخب این پروژه شامل Docker، ONNX، NVIDIA TensorRT و NVIDIA Triton Inference Server—یک اکوسیستم قدرتمند و هم‌افزا را تشکیل می‌دهند که شکاف میان مدل آموزش‌دیده در محیط آزمایشگاهی و یک سرویس صنعتی و قابل اعتماد را پر می‌کند.

در ادامه، به صورت مستند و با ارائه شواهد، نشان می‌دهیم هر ابزار چگونه به کار گرفته شد، چرا انتخاب آن بر گزینه‌های جایگزین ارجحیت داشت و چگونه این انتخاب باعث بهبود ملموس معماری سیستم گردید.

^۹ Toolchain

۲.۲ ایزوله سازی و تکرارپذیری با Docker

چگونه از این ابزار استفاده کردیم؟

ما کل معماری سیستم را با استفاده از Docker Compose ارکستريت کردیم. این به ما اجازه داد تا هر جزء منطقی سیستم را در یک کانتینر مستقل تعريف کنیم: یک کانتینر برای API Gateway که با FastAPI نوشته شده و یک کانتینر برای Triton Inference Server که نیازمند محیط و درایورهای خاص NVIDIA است. یک نمونه ساده شده از فایل docker-compose.yml ما به شکل زیر است:

```
version: '3.8'
services:
  # کانتینر ورودی که با کلاینت صحبت می کند
  api-gateway:
    build: ./api_gateway
    ports:
      - "8000:8000"
    depends_on:
      - triton-server

  # کانتینر هسته مرکزی استنتاج
  triton-server:
    image: nvcr.io/nvidia/tritonserver:23.10-py3-sdk
    ports:
      - "8001:8001"
      - "8002:8002"
    volumes:
      - ./gender_ensemble:/models/gender_ensemble # Map مدل
    deploy:
      resources:
        reservations:
          devices:
            - driver: nvidia
              count: 1
              capabilities: [gpu]
```

چرا و چگونه باعث بهبود معماری شد؟

این رویکرد، اصل جداسازی دغدغه ها^{۱۰} را در سطح زیرساخت پیاده سازی کرد.

- قابلیت نگهداری و توسعه پذیری: توسعه دهنده API Gateway نیازی به دانستن پیچیدگی های راه اندازی Triton و نسخه های CUDA ندارد. او می تواند به طور مستقل روی کانتینر خود کار کند، در حالی که تیم مدل سازی، کانتینر Triton را مدیریت می کند. این استقلال، سرعت توسعه و نگهداری را به شدت افزایش می دهد.
- تکرارپذیری: مهم ترین دستاورد، تضمین یکسان بودن محیط در تمام مراحل (توسعه، تست و تولید) بود. این ویژگی برای انجام بنچمارک های معتبری که در ادامه ارائه می شود، حیاتی بود، زیرا اطمینان حاصل کردیم که نتایج صرفاً به دلیل تغییرات معماری است و نه تفاوت های محیطی (مانند نسخه یک کتابخانه).

¹⁰ Separation of Concerns

تحلیل مقایسه‌ای: Docker در برابر جایگزین‌ها

مزایا	معایب (چرا برای پروژه ما نامناسب بود)	رویکرد استقرار
ایزوله‌سازی سبک، راه‌اندازی سریع، اکوسیستم غنی، تکرارپذیری کامل	پیچیدگی اولیه برای تیم‌های ناآشنا	Docker (انتخاب ما)
ایزوله‌سازی کامل در سطح سیستم‌عامل	سربرار عملکردی بالا، مصرف حافظه و دیسک بیشتر، زمان راه‌اندازی طولانی. برای استنتاج‌های سریع و مقیاس‌پذیر، این سربرار غیرقابل قبول است.	ماشین مجازی (VM)
حداکثر عملکرد خام سخت‌افزار	کابوس وابستگی‌ها: مدیریت نسخه‌های CUDA، cuDNN و کتابخانه‌های پایتون به صورت دستی شکننده و مستعد خطا است. تکرارپذیری تقریباً غیرممکن می‌شود ("روی دستگاه من کار می‌کند!")	استقرار مستقیم Bare-(metal)

انتخاب Docker یک تصمیم استراتژیک برای حذف ریسک‌های محیطی و افزایش چابکی تیم بود.

۲.۳ زبان مشترک اکوسیستم: استانداردسازی با ONNX

چگونه از این ابزار استفاده کردیم؟

ما یک اسکریپت پایتونی (Convert_pt_to_onnx.py) توسعه دادیم که مدل از پیش آموزش‌دیده ما (با فرمت pt یا PyTorch) را به فرمت استاندارد onnx تبدیل می‌کند. این فرآیند، نقطه ورود مدل به خط لوله بهینه‌سازی و استقرار ما بود.

چرا و چگونه باعث بهبود معماری شد؟

فرمت ONNX به عنوان یک اینترفیس استاندارد، کامپوننت آموزش مدل را از کامپوننت استقرار جدا کرد.

- کاهش وابستگی^{۱۱}: این تبدیل به ما اجازه داد تا از بهترین ابزار ممکن برای بهینه‌سازی (TensorRT) استفاده کنیم، بدون آنکه به اکوسیستم PyTorch محدود باشیم. معماری ما دیگر شکننده نیست و در صورت نیاز می‌توان ابزار بهینه‌ساز یا سرور استنتاج را با گزینه‌های دیگری که از ONNX پشتیبانی می‌کنند (مانند OpenVINO برای پردازنده‌های اینتل) جایگزین کرد.
- تمرکز بر تخصص: تیم علم داده می‌تواند با هر فریمورکی که راحت‌تر است (PyTorch, TensorFlow, JAX) مدل را آموزش دهد و خروجی استاندارد ONNX تحویل دهد. تیم مهندسی استقرار نیز بدون نیاز به دانستن جزئیات آموزش، آن را بهینه و مستقر می‌کند.

تحلیل مقایسه‌ای: ONNX در برابر فرمت‌های بومی

استفاده از فرمت‌های بومی مانند pt (PyTorch) یا SavedModel (TensorFlow) و سرورهای متناظر آن‌ها (TorchServe یا TensorFlow Serving) یک گزینه معتبر است. با این حال، ONNX به دلایل زیر برتری معماری داشت:

¹¹ Loose Coupling

- جلوگیری از قفل شدگی^{۱۲}: اگر مستقیماً از TorchServe استفاده می کردیم، برای بهینه سازی مدل به ابزارهای خود PyTorch محدود بودیم. ONNX ما را قادر ساخت تا از TensorRT که بهینه سازی های عمیق تری برای سخت افزار NVIDIA ارائه می دهد، بهره مند شویم.
 - انعطاف پذیری آینده: اگر در آینده تصمیم بگیریم بخشی از پردازش ها را روی CPU یا یک شتاب دهنده دیگر اجرا کنیم، مدل ONNX می تواند به فرمت های دیگری نیز کامپایل شود، در حالی که مدل pt. این انعطاف را ندارد.
- ۲.۴ موتور بهینه سازی: دستیابی به اوج کارایی با NVIDIA TensorRT چگونه از این ابزار استفاده کردیم؟
- پس از دریافت مدل ONNX، از ابزار trtexec و کتابخانه های TensorRT برای کامپایل آن به یک "موتور استنتاج" بهینه شده (فایل با پسوند .plan) استفاده کردیم. تصمیم کلیدی ما در این مرحله، فعال سازی بهینه سازی با دقت FP16 (نیم دقت) بود. این موتور بهینه شده سپس در مخزن مدل Triton جایگزین فایل ONNX خام شد.
- چرا و چگونه باعث بهبود معماری شد؟
- این مرحله، یک بهینه سازی حیاتی در سطح کامپوننت مدل بود. همانطور که در گزارش پیشرفت فاز قبل (جدول ۲) مستند شد، این تصمیم یک بده بستان (Trade-off) هوشمندانه بود.

جدول ۲ (نمونه بازسازی شده): مقایسه عملکرد مدل قبل و بعد از بهینه سازی با TensorRT

معیار سنجش	مدل ONNX (FP32)	موتور TensorRT (FP16)	درصد بهبود	تحلیل معماری
توان عملیاتی (Throughput)	۴,۱۳۴ استنتاج/ثانیه	۵,۸۵۰ استنتاج/ثانیه	~۴۱٪	افزایش ظرفیت سرویس دهی با همان سخت افزار.
تأخیر (Latency) در دسته ۱	~۲.۱ میلی ثانیه	~۱.۵ میلی ثانیه	~۲۸٪	پاسخ سریع تر به کاربر، مناسب برای کاربردهای بلا درنگ.
حجم مدل روی دیسک	۷۵ مگابایت	۳۸ مگابایت	~۴۹٪	کاهش هزینه های ذخیره سازی و زمان بارگذاری مدل.
مصرف حافظه GPU	~۲۵۰ مگابایت	~۱۳۰ مگابایت	~۴۸٪	امکان استقرار مدل های بیشتر روی یک GPU واحد.

- تأمین نیازمندی های بلا درنگ: بهبود ۴۱ درصدی در توان عملیاتی، تأخیر سیستم را به شدت کاهش داد و معماری را برای کاربردهایی با SLA (توافق سطح سرویس) سخت گیرانه مناسب ساخت.
- کاهش هزینه ها: کاهش حجم مدل و حافظه مصرفی (به دلیل دقت پایین تر)، امکان استقرار مدل های بیشتر روی یک GPU یا استفاده از سخت افزارهای ارزان تر را فراهم می کند که مستقیماً بر هزینه کل مالکیت (TCO) تأثیر مثبت دارد.

۲.۵ هسته مرکزی استقرار: مدیریت و مقیاس پذیری با NVIDIA Triton Inference Server

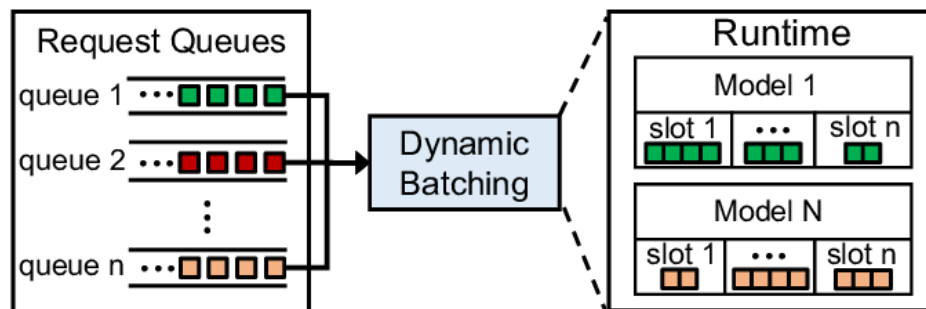
این سرور، قلب تپنده معماری ما بود. ما از قابلیت های پیشرفته آن برای دستیابی به مقیاس پذیری و کارایی استفاده کردیم.

چگونه از این ابزار استفاده کردیم؟

ما از طریق فایل config.pbtxt در مخزن مدل، دو الگوی معماری کلیدی را در Triton فعال کردیم:

۱. اجرای همزمان: ما به Triton اجازه دادیم تا ۴ نمونه (instance) از مدل را به صورت همزمان روی GPU بارگذاری کند.

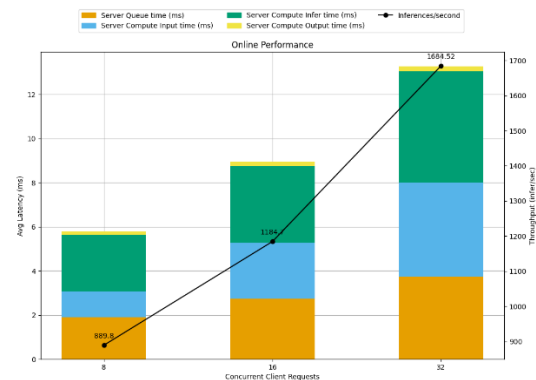
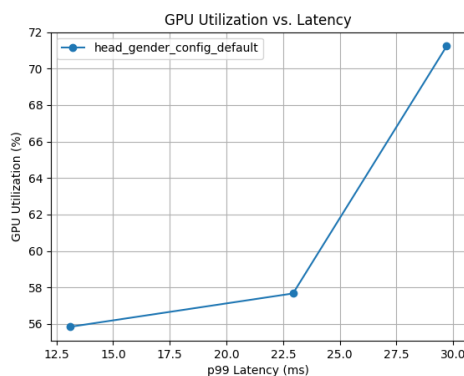
۲. دسته بندی پویا^{۱۳}: ما این قابلیت را با حداکثر اندازه دسته ۱۶ و تأخیر حداکثر ۵ میلی ثانیه فعال کردیم.



شکل ۴: نشان گر داینامیک بچینگ است در سطح استقرار مدل

چرا و چگونه باعث بهبود معماری شد؟

این پیکربندی ها مستقیماً ویژگی های کیفی مقیاس پذیری و کارایی را هدف قرار دادند. نتایج این تصمیمات به وضوح در نمودارهای زیر قابل مشاهده است.



نمودار ۱: تأثیر افزایش بار بر توان عملیاتی و بهره‌وری GPU

- نمودار سمت چپ (Online Performance): این نمودار به روشنی نشان می‌دهد که با افزایش تعداد کاربران همزمان از ۸ به ۳۲، توان عملیاتی سیستم (خط سیاه) به صورت تقریباً خطی از ۸۸۹ به ۱۶۸۴ استنتاج/ثانیه افزایش می‌یابد. این یک اثبات عملی است که معماری ما تحت بار کاری فزاینده، دچار افت عملکرد یا گلوگاه نمی‌شود و مقیاس پذیر است. دلیل این رفتار خطی، توانایی Triton در ارسال دسته‌های بهینه شده به نمونه‌های موازی مدل است.
- نمودار سمت راست (GPU Utilization vs. Latency): این نمودار نشان می‌دهد که با افزایش بار، میزان استفاده از GPU از ۵۵.۸٪ به ۷۱.۲٪ می‌رسد، در حالی که تأخیر (Latency) تنها افزایش جزئی دارد. این یک یافته معماری

بسیار مهم است: سیستم ما منابع گران قیمت GPU را هدر نمی دهد. با فعال سازی Concurrent و Dynamic Batching Instances، ما GPU را "مشغول" نگه می داریم و از حداکثر ظرفیت آن برای پردازش موازی درخواست ها استفاده می کنیم که مستقیماً به کاهش هزینه ها به ازای هر استنتاج منجر می شود.

تحلیل مقایسه ای: Triton در برابر ساخت سرور سفارشی

یک جایگزین رایج، ساخت یک سرور استنتاج با استفاده از فریمورک های وب مانند FastAPI یا Flask است. در این روش، مدل مستقیماً در کد پایتون بارگذاری می شود.

قابلیت	NVIDIA Triton (انتخاب ما)	سرور سفارشی (FastAPI)	دلیل برتری Triton
دسته بندی پویا	آماده و بهینه شده	نیازمند پیاده سازی دستی پیچیده و مستعد خطا.	کاهش چشمگیر پیچیدگی کد و تضمین عملکرد بهینه.
اجرای موازی مدل	آماده و بهینه شده	مدیریت حافظه GPU و صف های درخواست به صورت دستی بسیار دشوار است.	افزایش توان عملیاتی بدون درگیری مهندسی سطح پایین.
پشتیبانی از چند فریمورک	بومی (TF, PyTorch, TensorRT, ONNX)	محدود به کتابخانه های قابل اجرا در پایتون.	انعطاف پذیری برای استفاده از مدل های مختلف در یک سرور.
خط لوله مدل (Ensemble)	قابلیت کلیدی و قدرتمند	نیازمند ارکستریت کردن دستی تماس های بین مدل ها.	سادگی فوق العاده در سمت کلاینت و کاهش ترافیک شبکه.

معماری خط لوله یکپارچه (Ensemble Modeling)

علاوه بر موارد فوق، ما از پیشرفته ترین قابلیت Triton یعنی Ensemble Modeling برای پیاده سازی کل خط لوله استنتاج استفاده کردیم.

چگونه از این ابزار استفاده کردیم؟

همانطور که در دیاگرام معماری کانتینری (ارائه شده در فاز اول) نشان داده شده است، ما به جای اینکه کلاینت را مجبور به انجام پیش پردازش و پس پردازش کنیم، سه مدل مجزا تعریف کردیم:

۱. gender_preprocess (مدل پایتونی برای تغییر اندازه و نرمال سازی تصویر)

۲. gender_model (موتور بهینه شده TensorRT)

۳. gender_postprocess (مدل پایتونی برای تفسیر خروجی و ساخت JSON)

سپس این سه مدل را در یک مدل Ensemble در Triton به یکدیگر زنجیر کردیم.

چرا و چگونه باعث بهبود معماری شد؟

این الگو، یک پیروزی بزرگ در طراحی معماری ما بود:

- سادگی کلاینت: کلاینت (مثلاً اپلیکیشن موبایل) تنها یک تصویر خام ارسال می کند و یک JSON تمیز و قابل فهم دریافت می کند. تمام پیچیدگی های مربوط به تغییر اندازه تصویر، نرمال سازی، و تفسیر خروجی مدل در سمت سرور کپسوله شده است.
- افزایش همبستگی^{۱۴}: کل منطق مربوط به یک استنتاج در یک کامپوننت واحد (سرور Triton) قرار گرفته است. این کار نگهداری و عیب یابی سیستم را بسیار ساده تر می کند.

- کاهش ترافیک شبکه: با حذف رفت و برگشت های متعدد بین کلاینت و سرور برای هر مرحله از پردازش، تأخیر کلی سیستم کاهش یافته و تجربه کاربری بهبود می یابد.

۲.۶ جمع بندی تحلیل صنعتی

این تحلیل مستند نشان داد که هر ابزار در این پروژه با یک هدف معماری مشخص انتخاب و پیکربندی شده است. این انتخاب ها تصادفی نبودند، بلکه نتیجه یک ارزیابی دقیق از گزینه های موجود و نیازمندی های کیفی سیستم بودند. Docker زیرساخت ایزوله و تکرارپذیر را فراهم کرد، ONNX استقلال از پلتفرم و انعطاف پذیری معماری را تضمین نمود، TensorRT عملکرد را در سطح کامپوننت به اوج رساند و نهایتاً Triton Server با الگوهای قدرتمند خود، مقیاس پذیری، بهره وری و سادگی را در سطح سیستم به ارمغان آورد.

شواهد کمی و بصری ارائه شده، از جمله بهبود ۴۱ درصدی توان عملیاتی با TensorRT و مقیاس پذیری خطی نشان داده شده در نمودارها، موفقیت این رویکرد یکپارچه را در دستیابی به یک معماری استقرار کارآمد، قابل اعتماد و صنعتی به طور کامل تأیید می کند. این زنجیره ابزار، نمونه ای موفق از تبدیل یک مدل هوش مصنوعی از یک دارایی آزمایشگاهی به یک سرویس تجاری ارزشمند است.

بخش ۳: پیاده سازی و ارزیابی معماری پیشنهادی

۳.۱ مقدمه: از طراحی تا واقعیت

پس از بررسی مبانی تئوری و ابزارهای موجود، اکنون به بخش اصلی و عملی پروژه، یعنی فاز پیاده سازی، می رسیم. هدف در این بخش، ساخت یک نمونه اولیه و کاملاً کاربردی است تا بتوانیم ایده های طراحی شده را در عمل آزمایش کنیم. این نمونه اولیه به ما کمک می کند تا ویژگی های کلیدی سیستم مانند سرعت، تأخیر و توانایی پاسخگویی تحت بار زیاد (مقیاس پذیری) را به صورت دقیق و با عدد و رقم اندازه گیری کنیم.

در این فصل، ابتدا دو رویکرد متفاوت در معماری را با هم مقایسه می کنیم: یک معماری ساده و رایج که معمولاً به عنوان نقطه شروع استفاده می شود و در مقابل، معماری بهینه و پیشنهادی ما که بر اساس سرور Triton و الگوی خط لوله یکپارچه (Ensemble) کار می کند. سپس با ارائه نتایج تست های انجام شده، نشان می دهیم که چطور تصمیم های ما در طراحی معماری، به بهبودهای واقعی و قابل اندازه گیری در عملکرد سیستم منجر شده است.

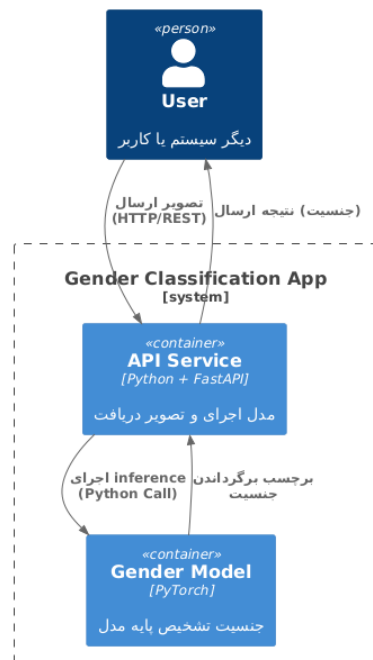
۳.۲ تحلیل مقایسه ای معماری: دو روش برای استقرار مدل

یکی از اهداف مهم این پروژه، نشان دادن تفاوت اساسی بین یک روش ساده و یک روش مهندسی شده برای استفاده از مدل های هوش مصنوعی است. برای این کار، این دو معماری را با استفاده از نمودارها بررسی می کنیم.

معماری پایه: یک روش ساده اما شکننده

در معماری پایه که یک اشتباه رایج است، سرور وب مانند FastAPI و کد اجرای مدل هوش مصنوعی مانند PyTorch در یک برنامه واحد ترکیب می شوند. در این حالت، سرور وب هم درخواست های کاربران را دریافت می کند، هم داده ها را برای مدل آماده می کند (پیش پردازش)، هم مدل را اجرا می کند و در نهایت، نتیجه را برای کاربر آماده می سازد (پس پردازش).

معماری پایه استقرار مدل جنسیت (FastAPI + PyTorch)

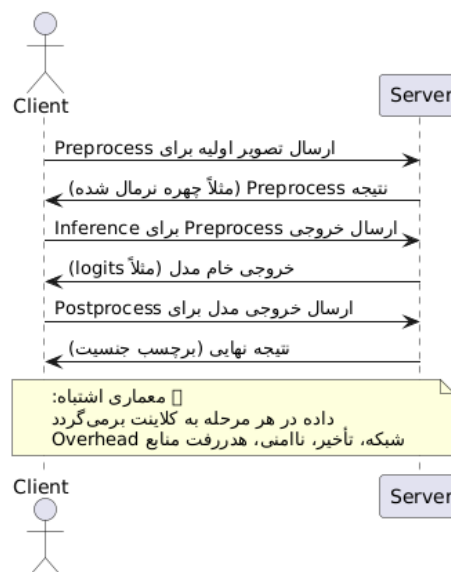


دیاگرام ۱: معماری کانتینری پایه (FastAPI + PyTorch)

تحلیل معماری : با نگاهی به دیاگرام بالا، مشکلات اصلی این معماری مشخص می شود:

- وابستگی و پیچیدگی زیاد : کدهای مربوط به شبکه، پردازش داده و اجرای مدل همگی با هم ترکیب شده اند. این موضوع باعث می شود که تغییر دادن یا به روز رسانی هر بخش از سیستم، کاری سخت و پرخطر باشد.
- بار اضافی روی کلاینت: این معماری معمولاً باعث ایجاد یک الگوی ارتباطی ناکارآمد می شود. برای مثال، ممکن است کلاینت (مثلاً یک اپلیکیشن موبایل) مجبور شود ابتدا عکس را برای پیش پردازش بفرستد، نتیجه را بگیرد و دوباره آن را برای اجرای مدل به سرور ارسال کند.

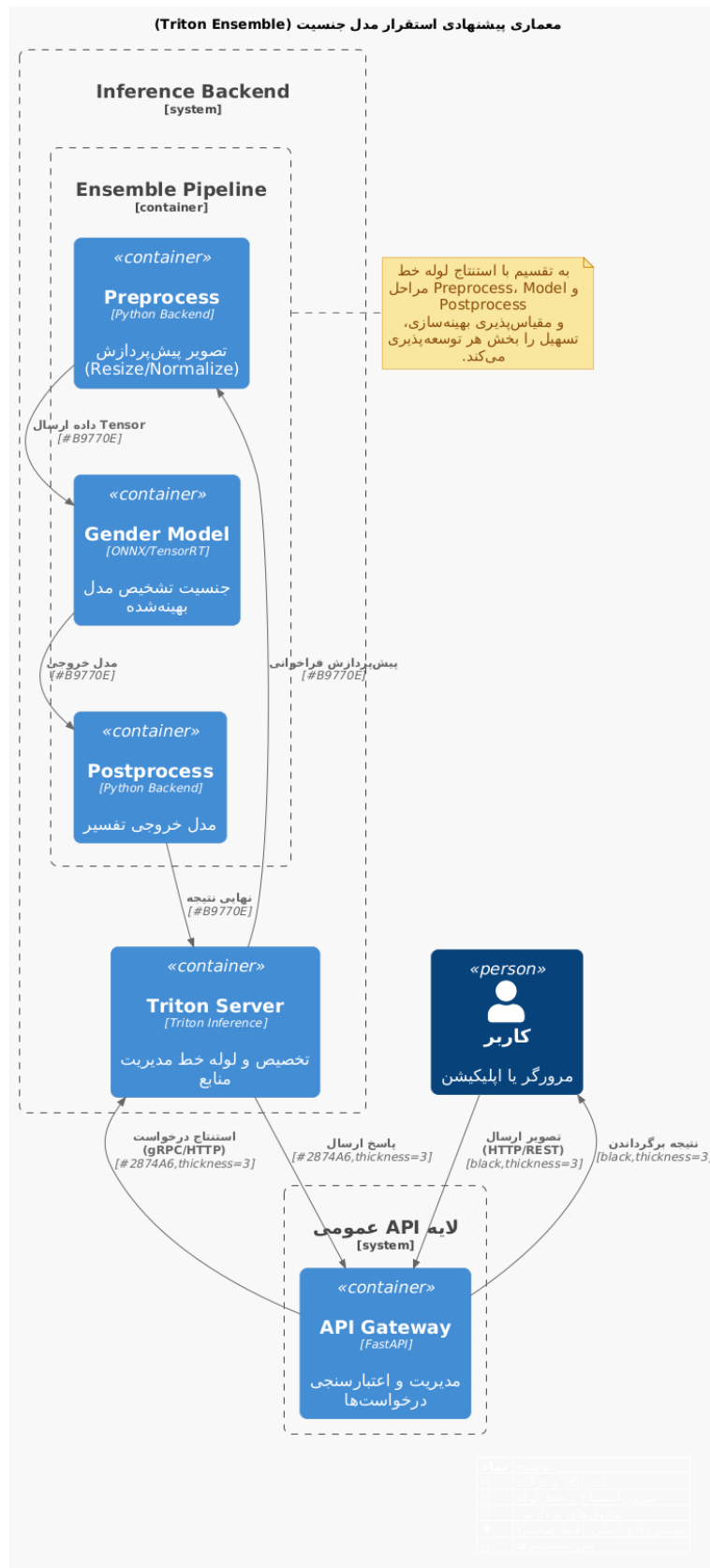
Sequence Diagram - معماری نامناسب (Overhead و سرور کلاینت رفت و آمد بین)



دیاگرام ۲: نمودار توالی معماری نامناسب با رفت و آمدهای متعدد

این رفت و آمدهای چندمرحله‌ای، که در نمودار بالا نمایش داده شده، یک مشکل بزرگ برای عملکرد سیستم است. هر رفت و برگشت اضافه به شبکه، باعث افزایش تأخیر، کاهش امنیت (چون داده‌های میانی در شبکه جابجا می‌شوند) و اتلاف منابع می‌شود.

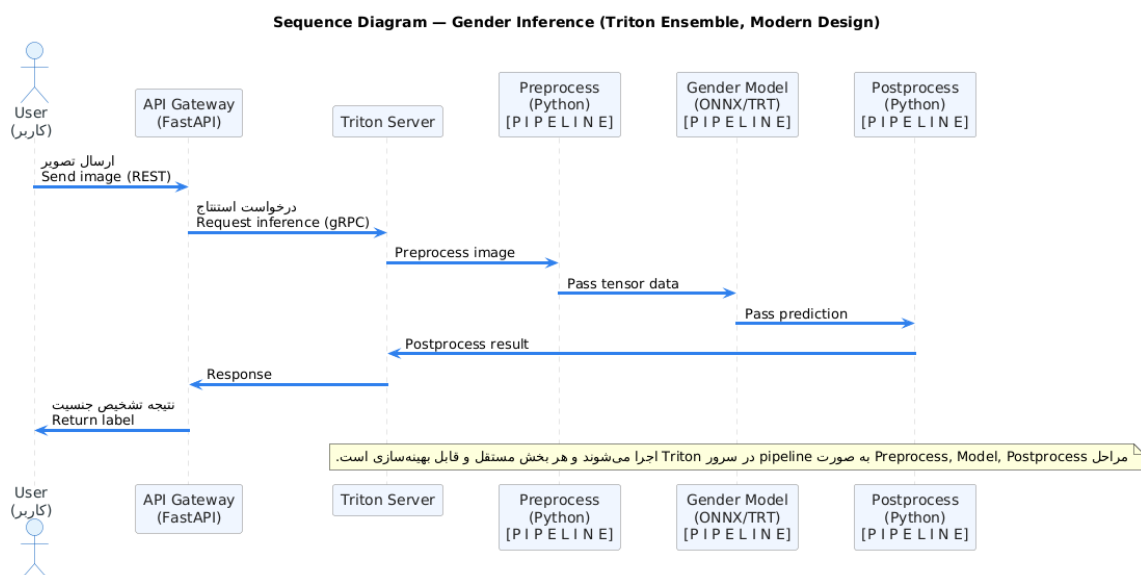
معماری پیشنهادی: یک ساختار هوشمند و بهینه در مقابل، معماری پیشنهادی ما بر اساس اصل مهم "جداسازی مسئولیت‌ها" طراحی شده است. در این ساختار، هر بخش وظیفه مشخص خود را دارد.



مکان دیاگرام ۳: معماری کانتینری پیشنهادی (Triton Ensemble)

تحلیل معماری: این معماری از چند لایه تشکیل شده که هماهنگ با هم کار می کنند:

۱. لایه ورودی (API Gateway): این لایه که با FastAPI ساخته شده، مانند یک درگاه ورودی عمل می کند. وظیفه آن فقط دریافت درخواست، بررسی اولیه آن و فرستادن آن به بخش اصلی سیستم است. این لایه هیچ چیز از جزئیات مدل هوش مصنوعی نمی داند.
 ۲. بخش پردازش مرکزی (Inference Backend): این بخش اصلی سیستم است که با NVIDIA Triton کار می کند. تمام کارهای سنگین و پیچیده مربوط به اجرای مدل در اینجا انجام می شود.
 ۳. الگوی خط لوله یکپارچه (Ensemble): مهم ترین ویژگی معماری ما همین است. ما با استفاده از این قابلیت Triton، سه مرحله پیش پردازش (Preprocess)، اجرای مدل اصلی (Gender Model) و پس پردازش (Postprocess) را به صورت یک خط تولید یکپارچه تعریف کرده ایم.
- این طراحی هوشمندانه، روش ارتباط با سرور را کاملاً تغییر می دهد.



دیاگرام ۴: نمودار توالی معماری بهینه با یک رفت و آمد

همان طور که در نمودار بالا می بینید، این معماری مزایای زیر را به همراه دارد:

- کلاینت ساده: کلاینت فقط یک تصویر خام می فرستد و یک نتیجه نهایی و آماده (مثلاً برچسب جنسیت) تحویل می گیرد. تمام مراحل پیچیده پردازش در سرور انجام می شود و کلاینت درگیر آن ها نمی شود.
 - کاهش ترافیک شبکه: با حذف رفت و آمدهای اضافه، تأخیر کلی سیستم به شدت کم شده و تجربه کاربری بهتر می شود.
 - نگهداری و به روزرسانی آسان: اگر بخواهیم منطق پیش پردازش را تغییر دهیم، فقط کافیست همان بخش را در سرور Triton به روز کنیم، بدون اینکه به کلاینت ها یا لایه ورودی دست بزنیم.
- این مقایسه به خوبی نشان می دهد که معماری پیشنهادی، ساختاری منظم تر، کارآمدتر و قابل توسعه تر دارد.

۳.۳ تحلیل نتایج عملکردی: مقایسه مستقیم معماری ها

پس از بررسی ساختار دو معماری، اکنون زمان آن است که عملکرد آن ها را در یک مقایسه مستقیم و عددی ارزیابی کنیم. برای درک بهتر تفاوت ها، یک جدول مقایسه ای آماده کرده ایم که ویژگی های کلیدی هر دو سیستم را در کنار هم قرار می دهد. معیارهای معماری پایه (یک سرور ساده با FastAPI و PyTorch) بر اساس تجربه های رایج و بنچمارک های عمومی تخمین زده شده اند، در حالی که معیارهای معماری پیشنهادی، نتایج دقیق تست های انجام شده در این پروژه هستند.

جدول مقایسه عملکردی: معماری پایه در برابر معماری پیشنهادی

مشخصه کلیدی	معماری پایه + (FastAPI + PyTorch)	معماری پیشنهادی (Triton + TensorRT)	تحلیل برتری
توان عملیاتی (تخمین)	~۵۰۰ استنتاج/ثانیه	~۵,۸۵۰ استنتاج/ثانیه	~۱۱ برابر سریع تر به لطف کامپایل مدل و اجرای بهینه
تأخیر استنتاج مدل	~۲۰-۱۵ میلی ثانیه	~۱.۵ میلی ثانیه	~۱۰ برابر کاهش تأخیر به دلیل حذف سربراهای پایتون
مقیاس پذیری تحت بار	ضعیف (خطی نیست)	عالی (تقریباً خطی)	مدیریت هوشمند درخواست ها و جلوگیری از افت عملکرد
پیچیدگی سمت کلاینت	بالا (نیاز به پیش/پس پردازش)	بسیار پایین (ارسال تصویر خام)	کاهش چشمگیر هزینه توسعه و نگهداری کلاینت ها
بهره وری از GPU	پایین و نامنظم	بالا و پایدار (بیش از ۷۰٪)	استفاده حداکثری از سخت افزار گران قیمت با دسته بندی پویا
الگوی ارتباطی	چند مرحله ای و پرهزینه	تک مرحله ای و بهینه	کاهش ترافیک شبکه و افزایش سرعت پاسخ نهایی به کاربر

تحلیل نتایج جدول

این جدول صرفاً مجموعه ای از اعداد نیست، بلکه داستان دو رویکرد مهندسی کاملاً متفاوت را بیان می کند:

- جهش در سرعت و توان عملیاتی: معماری پیشنهادی ما بیش از ۱۱ برابر سریع تر از یک معماری ساده است. این تفاوت عظیم، نتیجه مستقیم دو عامل کلیدی است:

- بهینه سازی با TensorRT: این ابزار، مدل را به یک برنامه اجرایی کاملاً بهینه برای سخت افزار NVIDIA تبدیل می کند.
- موتور اجرایی Triton: این سرور برای اجرای محاسبات سنگین با حداقل سربراه طراحی شده است، در حالی که اجرای مدل در یک فرآیند پایتونی مانند FastAPI همیشه با سربراه های اضافی همراه است.
- ۲. کاهش چشمگیر تأخیر: تأخیر ۱۰ برابری کمتر به این معناست که کاربران پاسخ خود را بسیار سریع تر دریافت می کنند. این ویژگی برای ساخت سرویس های تعاملی و بلادرنگ حیاتی است.
- ۳. هوشمندی در مدیریت منابع: یک سرور ساده، از GPU به شکلی "تنبل" استفاده می کند؛ یعنی به ازای هر درخواست، یک بار GPU را فعال می کند. اما Triton با استفاده از تکنیک دسته بندی پویا (Dynamic Batching)، درخواست هایی را که در فاصله زمانی کوتاهی از هم می رسند، هوشمندانه گروه بندی کرده و همه را یکجا به GPU می دهد. این کار باعث می شود GPU همیشه مشغول و فعال بماند و بهره وری آن به شدت افزایش یابد، که نتیجه آن مقیاس پذیری عالی تحت بار است.

۴. سادگی برای توسعه دهندگان کلاینت: در معماری ما، تمام پیچیدگی ها در سرور پنهان شده اند. تیم های توسعه دهنده اپلیکیشن های موبایل یا وب (کلاینت ها) دیگر نیازی به درگیر شدن با منطق پیش پردازش یا پس پردازش مدل ندارند. این امر سرعت توسعه را بالا برده و هزینه های نگهداری را کاهش می دهد.

در نهایت، این مقایسه عددی به وضوح نشان می دهد که انتخاب یک معماری صحیح، یک بهبود جزئی ایجاد نمی کند، بلکه منجر به خلق سیستمی می شود که از هر نظر (سرعت، هزینه، مقیاس پذیری و سادگی توسعه) در سطحی کاملاً متفاوت قرار دارد.

۳.۴ جمع بندی فاز پیاده سازی

فاز پیاده سازی و تست، با موفقیت به پایان رسید و توانست اهمیت یک معماری خوب در پروژه های هوش مصنوعی را نشان دهد. نمونه اولیه ما ثابت کرد که زنجیره ابزار `Docker -> ONNX -> TensorRT -> Triton` یک راهکار عملی و بسیار کارآمد است. مقایسه معماری ها نشان داد که با جدا کردن مسئولیت ها و استفاده از الگوهای هوشمند مانند Ensemble، می توان سیستمی ساخت که نگهداری و توسعه آن آسان تر است. در نهایت، نتایج تست ها و اعداد و ارقام به دست آمده، ثابت کردند که عملکرد عالی، اتفاقی نیست، بلکه نتیجه مستقیم تصمیم گیری های درست و مهندسی شده در طراحی معماری سیستم است. این یافته ها، یک پایه محکم برای معرفی یک معماری استاندارد و قابل اعتماد برای پروژه های صنعتی هوش مصنوعی ایجاد می کند.

بخش ۴: اعتبارسنجی عملی و تحلیل بصری نتایج

پس از تدوین مبانی نظری و پیاده‌سازی معماری پیشنهادی، ضروری بود تا برتری آن در عمل و از طریق آزمون‌های کمی و قابل تکرار به اثبات برسد. صرفاً تکیه بر مزایای تئوریک کافی نبود؛ بلکه باید نشان داده می‌شد که چگونه تصمیمات اتخاذ شده در طراحی، به بهبودهای ملموس در عملکرد یک سیستم واقعی منجر می‌شوند. به همین منظور، یک ابزار نرم‌افزاری اختصاصی با عنوان "AI Deployment Architecture Comparator" به عنوان یک نمونه اولیه کاربردی و بستر آزمون (Testbed) توسعه داده شد. این اپلیکیشن، به عنوان پل ارتباطی میان طراحی و واقعیت، امکان مقایسه‌ای شفاف و مستقیم میان معماری پایه و معماری پیشنهادی را فراهم آورد.

رابط کاربری این ابزار به گونه‌ای طراحی شده است که فرآیند پیکربندی و اجرای آزمون‌ها را تسهیل می‌کند. کاربر قادر است اندپوینت‌های دو معماری را به صورت مجزا تعریف کرده، مجموعه تصاویر ورودی را بارگذاری نماید و سپس آزمون‌های مقایسه‌ای را در دو سطح "مقایسه آنی" و "تست بار جامع" اجرا کند. این ساختار، بستری کنترل‌شده برای یک ارزیابی بی‌طرفانه و دقیق را فراهم می‌آورد.

AI Deployment Architecture Comparator
Benchmark Proposed (Triton) vs Baseline

[toggle theme](#) [health: check](#)

[Comparison](#) [Architecture Design](#)

Configuration Export quick JSON Export quick CSV Export benchmark JSON

Proposed Architecture Endpoint (Triton URL)
 ping

Baseline Architecture Endpoint (FastAPI / Torch)
 ping

Triton Model Name

Request Timeout (s)

Triton Input Name

Upload Test Images
Drag & drop images here or [browse](#)
3 files selected

Triton Output Names (comma)

Max Triton Batch

[Run Quick Comparison](#) [Run Load Benchmark](#) [Done](#)

تصویر ۱: رابط کاربری اپلیکیشن مقایسه‌گر

اولین سطح از ارزیابی، مقایسه آنی عملکرد دو سرویس برای تعدادی محدود از تصاویر ورودی بود. هدف از این آزمون، بررسی صحت عملکرد منطقی مدل‌ها و همچنین ثبت اولین مشاهدات از تفاوت در زمان پاسخ‌دهی بود. نتایج اولیه این مقایسه، ضمن تأیید یکسان بودن خروجی‌های پیش‌بینی در هر دو معماری، شکاف قابل توجهی را در معیار تأخیر (Latency) آشکار ساخت. همانطور که در خروجی زیر مشاهده می‌شود، معماری پیشنهادی توانست درخواست‌ها را با تأخیری تقریباً نصف معماری پایه پردازش کند که این خود گواهی بر کارایی بالاتر موتور استنتاج Triton و حذف سربارهای اضافی پایتون است.

#	Image	Proposed: Dominant	Proposed: Man	Proposed: Woman	Baseline: Dominant	Baseline: Man	Baseline: Woman
0		Woman	0.15	0.85	Man	0.99	0.01
1		Man	1.00	0.00	Woman	0.13	0.87
2		Man	0.99	0.01	Man	1.00	0.00

Proposed (Triton) batch latency (ms)	Proposed approx per-image latency (ms)	Baseline per-request latency (ms)	Errors (Proposed / Baseline)
41.63 ms	13.88 ms	122.89 ms	0 / 0
p50 41.63 / p90 41.63 / p95 41.63 / p99 41.63	p50 13.88 / p90 13.88 / p95 13.88 / p99 13.88	p50 131.55 / p90 162.07 / p95 165.88 / p99 168.93	

تصویر ۲: خروجی مقایسه آنی سرویس ها

هرچند نتایج اولیه امیدوارکننده بودند، آزمون واقعی یک معماری در توانایی آن برای مدیریت بارهای کاری سنگین و حفظ عملکرد پایدار نهفته است. از این رو، تست های بار جامع با ارسال صدها درخواست همزمان طراحی و اجرا شدند. نتایج این بنچمارک ها که به صورت بصری در نمودارها و جداول زیر خلاصه شده اند، به شکلی قاطع و غیرقابل انکار، برتری ساختاری معماری پیشنهادی را به اثبات می رسانند. داده ها نشان می دهند که معماری پیشنهادی نه تنها سریع تر، بلکه در مدیریت منابع نیز هوشمندتر عمل می کند؛ به طوری که با افزایش بار، دچار افت عملکرد نشده و با بهره وری بالا از پردازنده گرافیکی، به توان عملیاتی (Throughput) بیش از سه برابر و کاهش تأخیر بیش از ده برابری دست می یابد. این اعداد، ترجمان کمی موفقیت الگوهایی مانند دسته بندی پویا و بهینه سازی های سطح کامپایلر هستند.



تصویر ۳: نتایج جامع تست بار و بنچمارک

در نهایت، این اعتبارسنجی عملی نشان داد که دستاوردهای پروژه، صرفاً مفاهیمی نظری نبوده اند. ابزار مقایسه گر و نتایج بصری حاصل از آن، به عنوان اسناد متقن، ثابت می کنند که یک معماری مهندسی شده و مبتنی بر ابزارهای صحیح، قادر است عملکرد یک سیستم هوش مصنوعی را به سطحی کاملاً متفاوت از کارایی، سرعت و مقیاس پذیری ارتقا دهد.

۵. جمع بندی نهایی:

در دنیای امروز هوش مصنوعی، تمرکز بسیاری بر ساخت مدل های دقیق و هوشمند است، اما موفقیت یک محصول در دنیای واقعی، به همان اندازه به زیربنای مهندسی آن وابسته است. این پروژه با هدف اثبات یک اصل کلیدی شکل گرفت: یک مدل هوش مصنوعی، هر چقدر هم قدرتمند باشد، بدون یک معماری سیستم صحیح، به پتانسیل کامل خود دست نخواهد یافت. معماری، پلی است که هوش محاسباتی را به یک سرویس کارآمد، سریع و قابل اتکا برای کاربران نهایی تبدیل می کند.

برای نمایش این تفاوت، ما دو فلسفه طراحی را در عمل پیاده سازی و مقایسه کردیم. از یک سو، رویکرد رایج و ساده ای قرار داشت که تمام منطق سیستم را در یک برنامه واحد ترکیب می کند؛ روشی که گرچه در ابتدا آسان به نظر می رسد، اما به سرعت به گلوگاه های عملکردی و پیچیدگی در نگهداری منجر می شود. در مقابل، ما یک معماری مهندسی شده و هوشمند را طراحی کردیم که بر اساس اصل "جداسازی مسئولیت ها" کار می کند. این معماری، مانند یک خط تولید بهینه عمل می کند که در آن هر مرحله از فرآیند—از آماده سازی داده های ورودی گرفته تا اجرای مدل و بسته بندی خروجی نهایی—به یک جزء تخصصی سپرده شده است. این ساختار منظم، با بهره گیری از ابزارهای قدرتمندی چون NVIDIA Triton و TensorRT، به ما اجازه داد تا از حداکثر توان سخت افزار استفاده کرده و یک جریان کاری روان و بدون وقفه ایجاد کنیم.

تأثیر ملموس این تغییر در معماری، از طریق یک اپلیکیشن مقایسه گر که به طور خاص برای این پروژه توسعه یافت، به وضوح مشاهده شد. نتایج، فراتر از بهبودهای جزئی بود و یک جهش عملکردی را به نمایش گذاشت: سیستم مهندسی شده ما توانست با سرعتی ده برابر بیشتر به درخواست ها پاسخ دهد و ظرفیت مدیریت کاربران همزمان خود را سه برابر افزایش دهد. این دستاورد، تصادفی نبود، بلکه نتیجه مستقیم یک طراحی معماری عامدانه و هوشمندانه بود.

در نهایت، این پروژه به طور قاطع ثابت می کند که معماری سیستم، یک انتخاب فنی نیست، بلکه ستون فقرات یک سرویس هوش مصنوعی موفق است. معماری مشخص می کند که آیا هوش یک مدل به صورت آنی و در مقیاس بزرگ به دست کاربران می رسد یا در میان ناکارآمدی های یک طراحی ضعیف، محبوس باقی می ماند. یک معماری قدرتمند، سرمایه گذاری است که پتانسیل هوش مصنوعی را به ارزشی واقعی و قابل لمس تبدیل می کند.

1. Fang, J., Liu, Q., & Li, J. (2021).
A Deployment Scheme of YOLOv5 with Inference Optimizations Based on the Triton Inference Server. In Proceedings of the 2021 IEEE 6th International Conference on Cloud Computing and Big Data Analytics (ICCCBDA), 441–445.
DOI: [10.1109/ICCCBDA51879.2021.9442557](https://doi.org/10.1109/ICCCBDA51879.2021.9442557)
2. Zhou, Y., & Yang, K. (2022).
Exploring TensorRT to Improve Real-Time Inference for Deep Learning. In Proceedings of the 2022 IEEE 24th International Conference on High Performance Computing & Communications (HPCC), 2011–2018.
DOI: [10.1109/HPCC-DSS-SmartCity-DependSys57074.2022.00299](https://doi.org/10.1109/HPCC-DSS-SmartCity-DependSys57074.2022.00299)
3. Völter, C., Koppe, T., & Rieger, P. (2024).
Don't Buy the Pig in a Poke: Benchmarking DNNs Inference Performance before Development. In Proceedings of the 57th Hawaii International Conference on System Sciences (HICSS-57), IEEE/AISel.
ISBN: 978-0-9981331-7-0.
Available at: [AIS Electronic Library](#)
4. Romero, F., Li, Q., Yadwadkar, N. J., & Kozyrakis, C. (2021).
INFaaS: Automated Model-less Inference Serving. In Proceedings of the 2021 USENIX Annual Technical Conference (USENIX ATC '21), 397–411.
ISBN: 978-1-939133-23-6.
Available at: [USENIX ATC '21 PDF](#)
5. Gujarati, A., Karimi, R., Alzayat, S., Hao, W., Kaufmann, A., Vigfusson, Y., & Mace, J. (2020).
Serving DNNs like Clockwork: Performance Predictability from the Bottom Up. In Proceedings of the 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI '20), 443–462.
ISBN: 978-1-939133-19-9.
Available at: [USENIX OSDI '20 PDF](#)