

Problema

Qual modelo escolher?

Trade-off principal:

- Velocidade **vs** Qualidade
- Recursos **vs** Detalhamento

Modelos Comparados:

Modelo	Parâmetros	Tamanho
Flan-T5-Small	80M	320 MB
Flan-T5-Base	250M	990 MB

Como Testamos

Executamos 20 testes em 5 categorias:

- Conhecimento geral (4 testes)
- Matemática (4 testes)
- Tradução (3 testes)
- Lógica (3 testes)
- Criatividade (3 testes)

Solução

Código para comparar modelos

```
from transformers import pipeline

# Carregar modelos do HuggingFace
model_small = pipeline("text2text-generation",
                      model="google/flan-t5-small")

model_base = pipeline("text2text-generation",
                      model="google/flan-t5-base")

# Testar e medir tempo
for pergunta in perguntas:
    inicio = time.time()
    resposta = model(pergunta)
    tempo = time.time() - inicio

    # Salvar métricas
    resultados.append({
        'modelo': nome,
        'tempo': tempo,
        'resposta': resposta
    })
```

O que medimos:

- Tempo de resposta (segundos)
- Velocidade (tokens/segundo)
- Tamanho das respostas
- Uso de memória

Tecnologias:

- HuggingFace Transformers - Carregar modelos
- PyTorch - Backend
- Pandas - Análise de dados
- Matplotlib - Gráficos

Execução

