

Tokens em Modelos de Linguagem e sua Importância para Agentes de IA

1. O que são tokens

Tokens são as menores unidades de texto que um modelo de linguagem consegue processar.

Um texto como:

```
"Aprender IA é incrível!"
```

Pode ser transformado em algo como:

```
["Aprender", " IA", " é", " incrível", "!"]
```

Ou até em pedaços menores, dependendo do modelo.

Importante:

O modelo não entende palavras diretamente — ele entende **tokens**.

2. Tipos de tokens

2.1 Tokens comuns (tokens de texto)

São os tokens que representam conteúdo normal:

- palavras
- partes de palavras
- pontuação

Exemplo:

```
"inteligência" → ["inteli", "gência"]
```

2.2 Tokens especiais

São tokens que **não representam texto comum**, mas sim **estrutura ou controle**.

Eles funcionam como instruções internas para o modelo.

3. Principais tipos de tokens especiais

3.1 BOS (Beginning Of Sequence)

Indica o início da entrada.

```
<bos>
```

Serve para avisar ao modelo:

“Aqui começa o conteúdo”

3.2 EOS (End Of Sequence)

Indica o fim da geração.

```
<eos>
```

Função:

- Diz ao modelo quando parar
- Evita geração infinita

Analogia:

É como o ponto final de uma frase, mas com poder de encerrar completamente a resposta.

3.3 PAD (Padding)

Usado para completar sequências até um tamanho fixo.

```
<pad>
```

Muito usado em treinamento.

3.4 Tokens de papel (role tokens)

Usados em chats:

```
<user><assistant><system>
```

Servem para indicar:

- quem está falando
 - contexto da conversa
-

3.5 Tokens delimitadores

Alguns modelos usam formatos próprios:

```
[INST] ... [/INST]
```

ou

```
<s> ... </s>
```

Eles delimitam blocos de instrução.

4. Como isso funciona em um modelo de chat

Por trás de uma conversa simples, o modelo recebe algo como:

```
<system> Você é um assistente útil </system><user> O que é IA? </user><assistant>
```

O modelo então completa:

```
IA é o campo da computação que... <eos>
```

5. Por que tokens especiais são importantes

5.1 Estrutura da conversa

Sem tokens especiais:

```
Usuário: Oi Assistente: Olá Usuário: Tudo bem?
```

Para o modelo isso vira um texto confuso.

Com tokens:

```
<user> Oi </user><assistant> Olá </assistant><user> Tudo bem? </user>
```

Agora há estrutura clara.

5.2 Controle da geração

O token `<eos>` permite:

- parar respostas automaticamente
 - evitar loops infinitos
 - melhorar performance
-

5.3 Definição de comportamento

O token `<system>` pode mudar completamente o modelo:

```
<system> Seja formal </system>
```

Isso altera o estilo da resposta.

5.4 Compatibilidade entre modelos

Cada modelo tem seu próprio formato.

Exemplo:

- Llama → usa `[INST]`
- OpenAI → usa roles (system/user/assistant)
- Outros → usam JSON interno

Por isso existem **chat templates**:

“ eles adaptam a entrada para o formato correto do modelo

6. Importância para agentes de IA

Agentes de IA dependem fortemente desses tokens porque:

6.1 Mantêm contexto

Permitem separar:

- instruções

- histórico
 - respostas
-

6.2 Controlam comportamento

O agente pode:

- mudar personalidade
- seguir regras
- executar tarefas específicas

Tudo via estrutura de tokens.

6.3 Evitam ambiguidade

Sem tokens especiais, o modelo pode:

- confundir quem está falando
 - responder de forma incoerente
-

6.4 Permitem automação

Ferramentas, memória e raciocínio estruturado dependem de:

- delimitação clara
 - início/fim de blocos
 - controle de geração
-

7. Insight avançado

Tokens especiais são uma forma de:

“programar o modelo usando texto”

Eles funcionam como uma camada de controle sem precisar alterar o código do modelo.

8. Conclusão

- Tokens são a base do funcionamento dos LLMs
- Tokens especiais organizam e controlam o comportamento
- O token EOS é essencial para indicar o fim da resposta
- Em agentes de IA, eles são fundamentais para:
 - contexto
 - estrutura
 - controle
 - previsibilidade

Revision #3

Created 18 March 2026 14:27:25 by Lucas Arruda

Updated 29 May 2026 18:52:28 by Lucas Arruda