## The Pivot Word

Naïve Proof of Inference for Autoregressive Models

Vitaly Kleban, Brussels, 2024 @vkleban Alice wants to run an <u>Llama-3-8B</u> model inference on a set of inputs but doesn't have enough computing power.

**Bob** has the necessary computing power and offered to help Alice by performing the inference for her.

Now, **Alice** wants to verify that **Bob** <u>actually used</u> the specified model for the inference.

Alice needs to spend much less time verifying the input rather than doing the inference.

**Bob** (Llama-3-8B, seed=42, temperature=1)

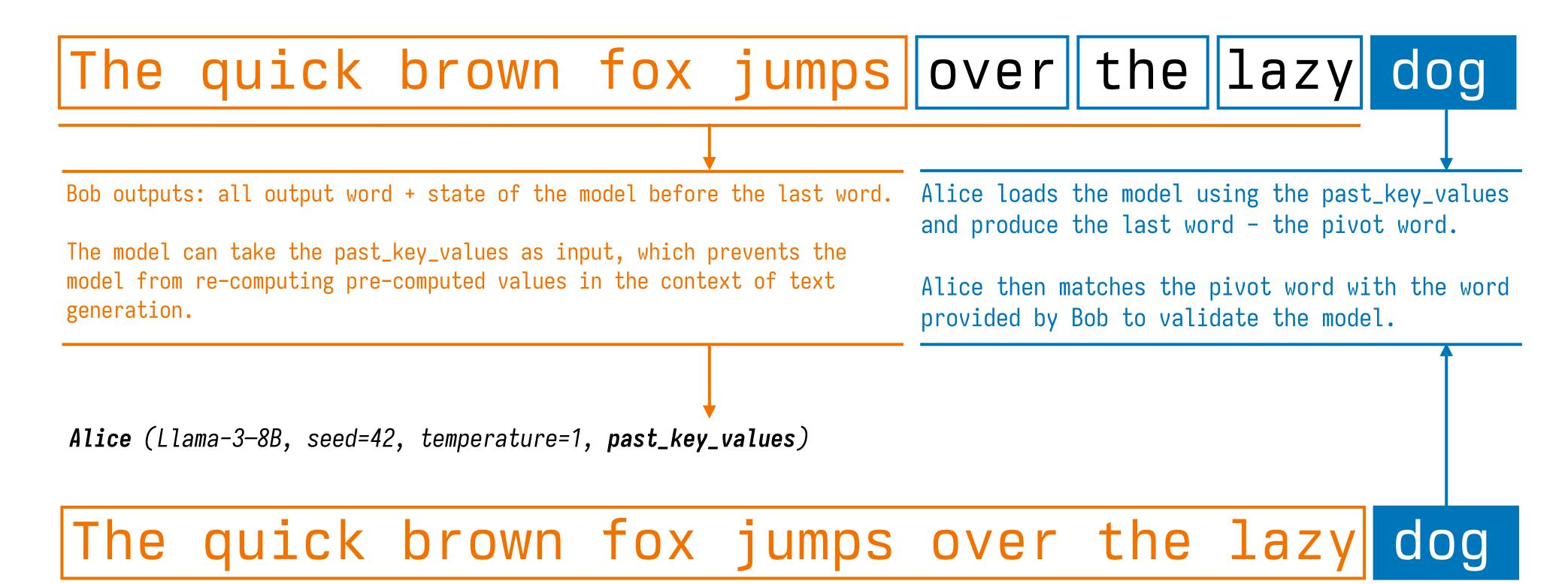
## The quick brown fox jumps over the lazy dog

dependent on the sequence of tokens that have been provided as input up to that point and the state of the model (previously computed key/value attention pairs).

The next token generated by the LLM is

Alice (Llama-3-8B, seed=42, temperature=1)

The quick brown fox jumps over the lazy dog



• The model can take the <code>past\_key\_values</code> (for PyTorch) or <code>past</code> (for TF) as input, which is the previously computed key/value attention pairs. Using this (<code>past\_key\_values</code> or <code>past</code>) value prevents the model from re-computing pre-computed values in the context of text generation. For PyTorch, see <code>past\_key\_values</code> argument of the <code>GPT2Model.forward()</code> method, or for TF the <code>past</code> argument of the <code>TFGPT2Model.call()</code> method for more information on its usage.

Bob processes all the words and shares the results with Alice.

The final word, the pivot, will determine if Alice trusts Bob.

Bob also sends the attention values recorded just before generating the pivot word.

Alice gets the output, including the pivot word from Bob.

She loads the model with Bob's attention values and runs a one-word inference to see if it matches the pivot word.

Chat with me at @vkleban!

Send your pitch to <u>vk@cyber.fund</u>

Check out "salty embeddings" from the same author <a href="https://github.com/vkleban/salty-embeddings">https://github.com/vkleban/salty-embeddings</a>