

Input "Welcome back ladies and"

Tokens Welcome back ladies and

Embeddings



Gradient-based Attribution Scores



Transformer LM

1. Calculate prediction



0.2% ... 90%

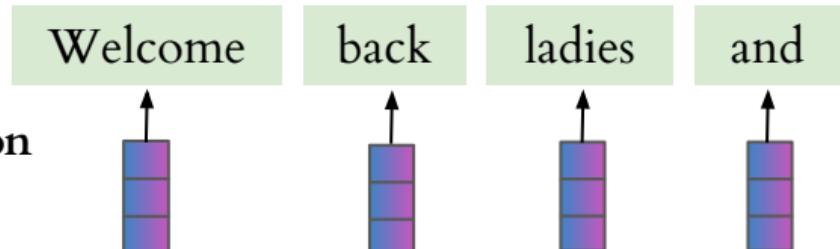
aadvark
gentlemen

2. Select a target token



90% gentlemen

4. Token Aggregation



3. Backpropagate gradients