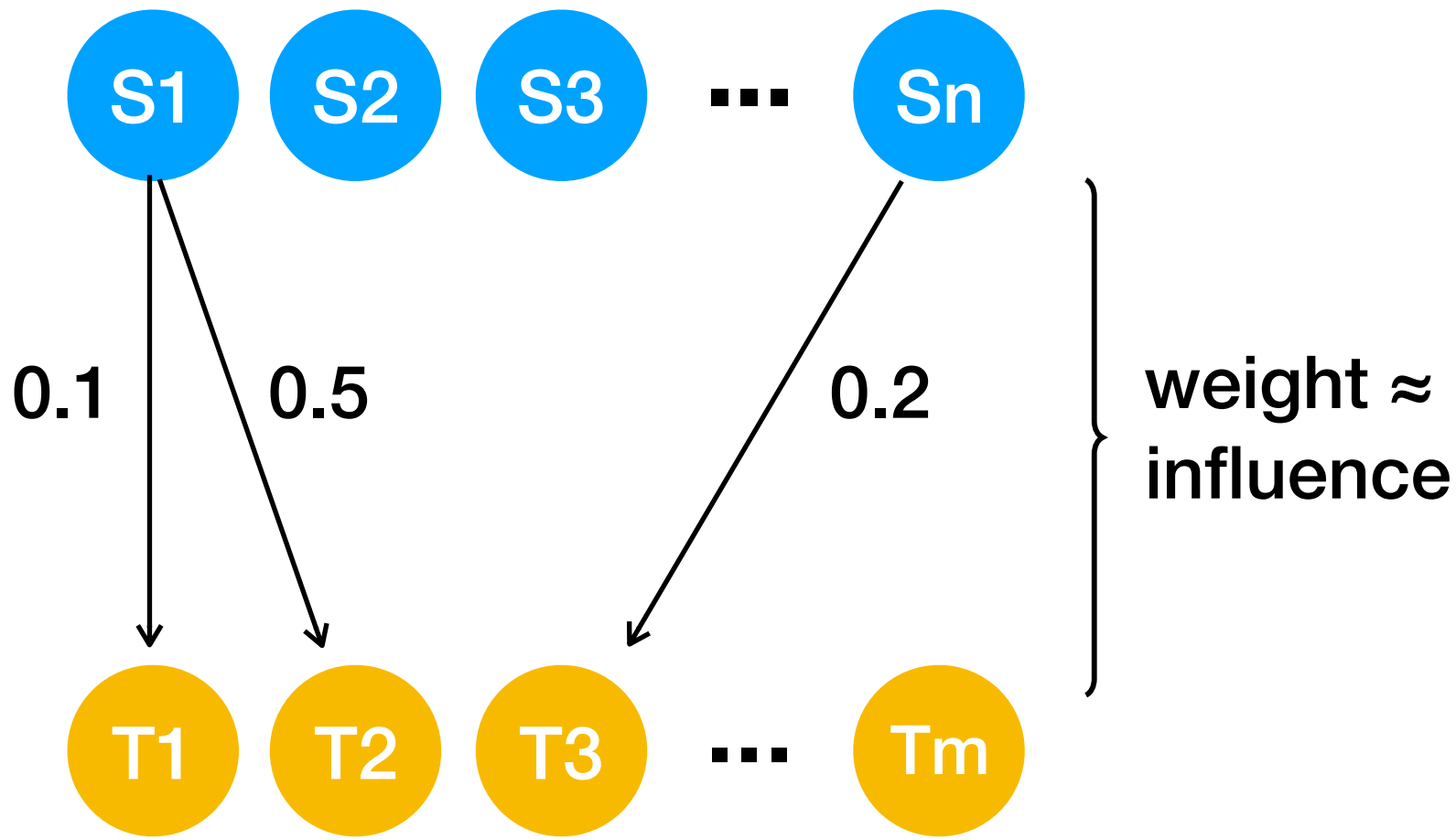




# Interpretability for Black-Box Seq2Seq Models

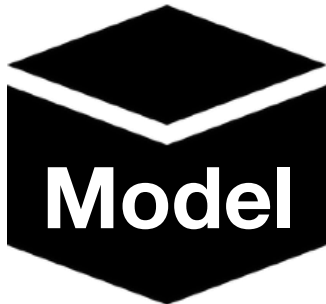


ANN + Jaakkola, ENR'17



**Input:**

$(S1, S2, \dots, S_n)$



**Output:**

$(T1, T2, \dots, T_n)$

**(locally)**



• Weighted bipartite graph summarizes model's behavior:



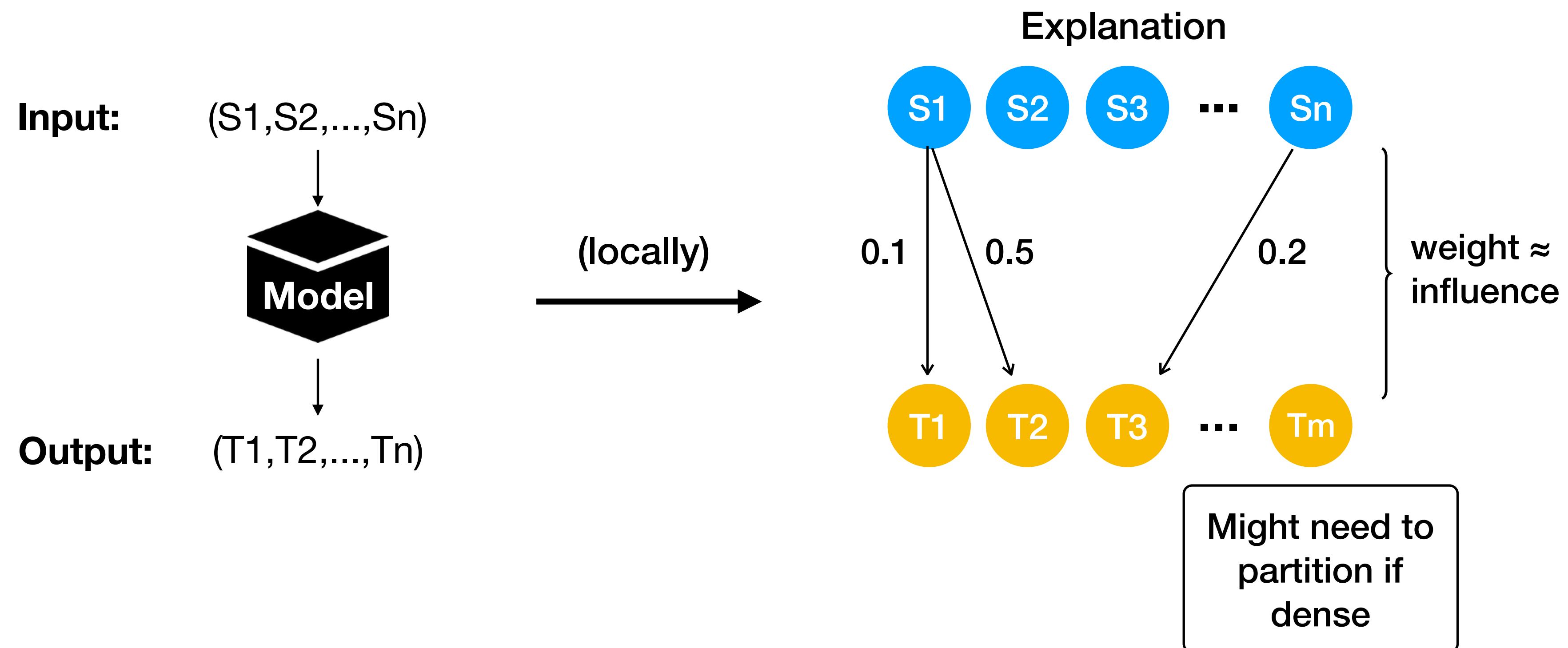
**Explanation**

**Might need to  
partition if  
dense**

# Interpretability for Black-Box Seq2Seq Models

AM+ Jaakkola, *EMNLP'17*

- Weighted bipartite graph summarizes model's local behavior:



# Interpretability for Black-Box Seq2Seq Models

AM+ Jaakkola, *EMNLP'17*

