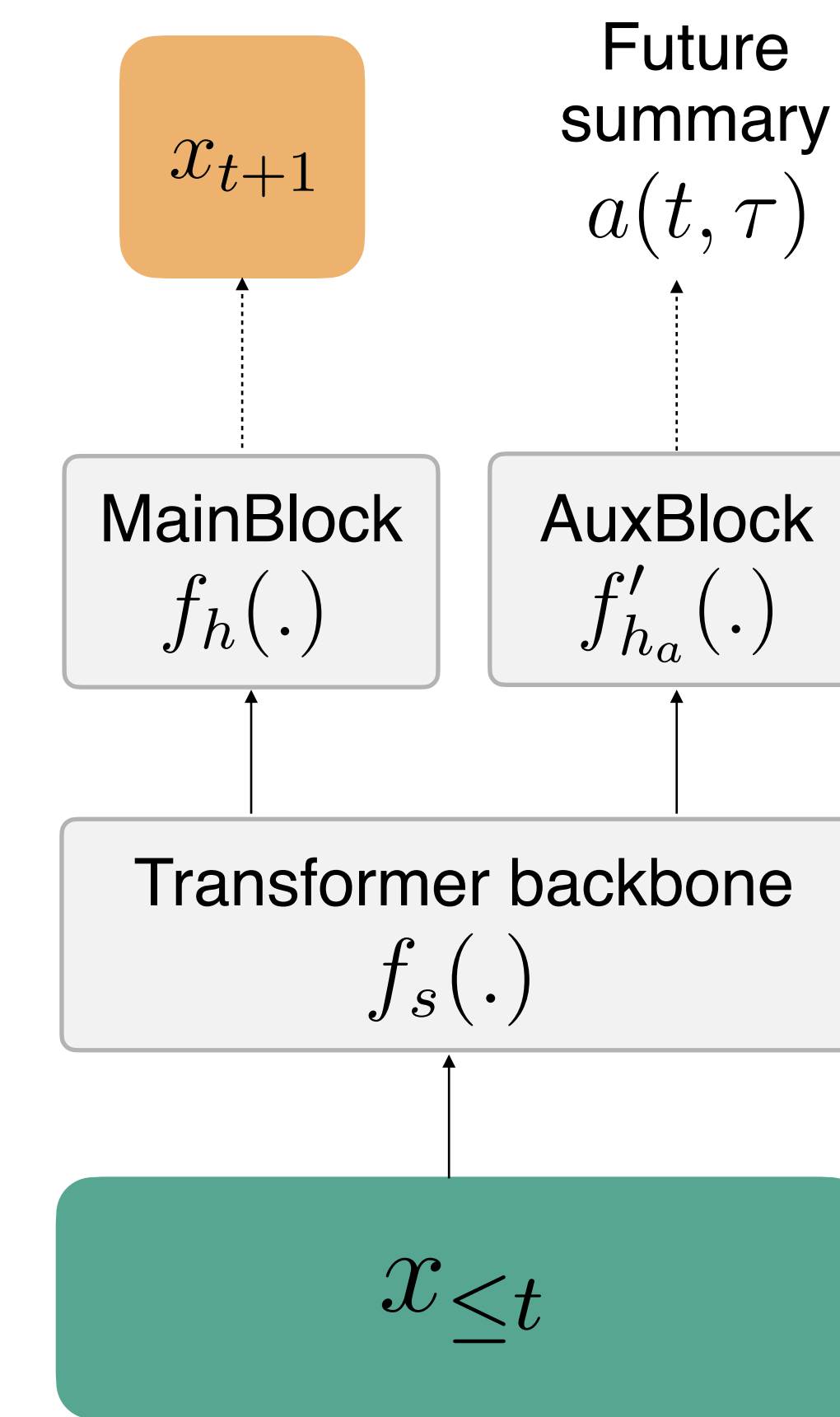


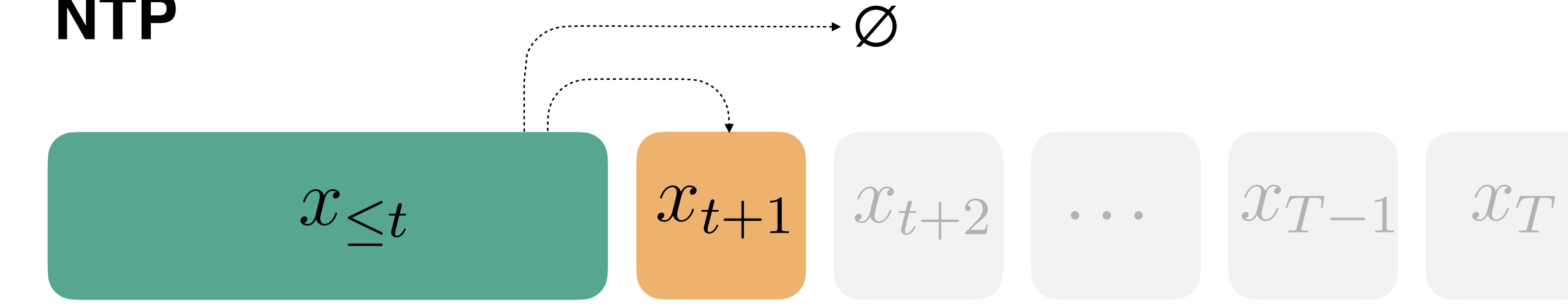


Contributions

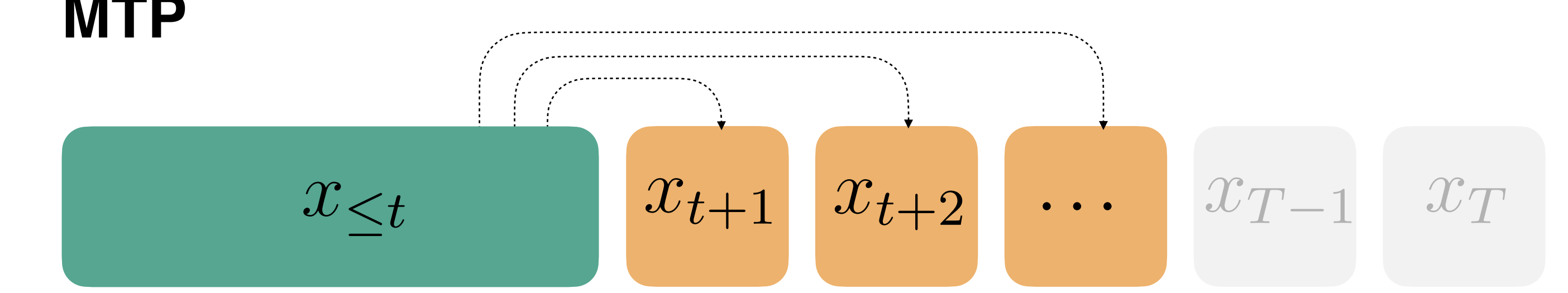
- Limitations of *Next-token prediction* (NTP) due to teacher forcing
 - Exposure bias & Shortcut learning
- *Multi-token prediction* (MTP): Predict short blocks of future tokens to reduce teacher forcing
 - One auxiliary head per future token is not scalable
- Proposal: *Future summary prediction* (FSP) to capture long-range dependencies
 - *Single auxiliary head* predicts a compact summary of future sequences
- Techniques for future summary prediction
 - *Hand-crafted summary*: Bag-of-words as auxiliary target
 - *Learned summary*: Embedding from reverse language model as auxiliary target



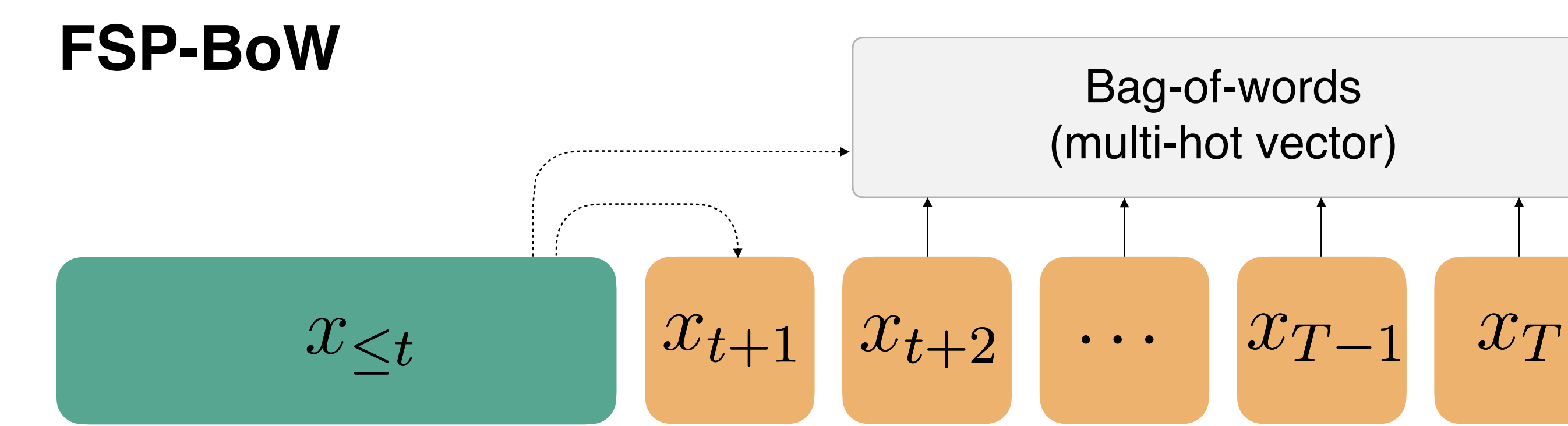
NTP



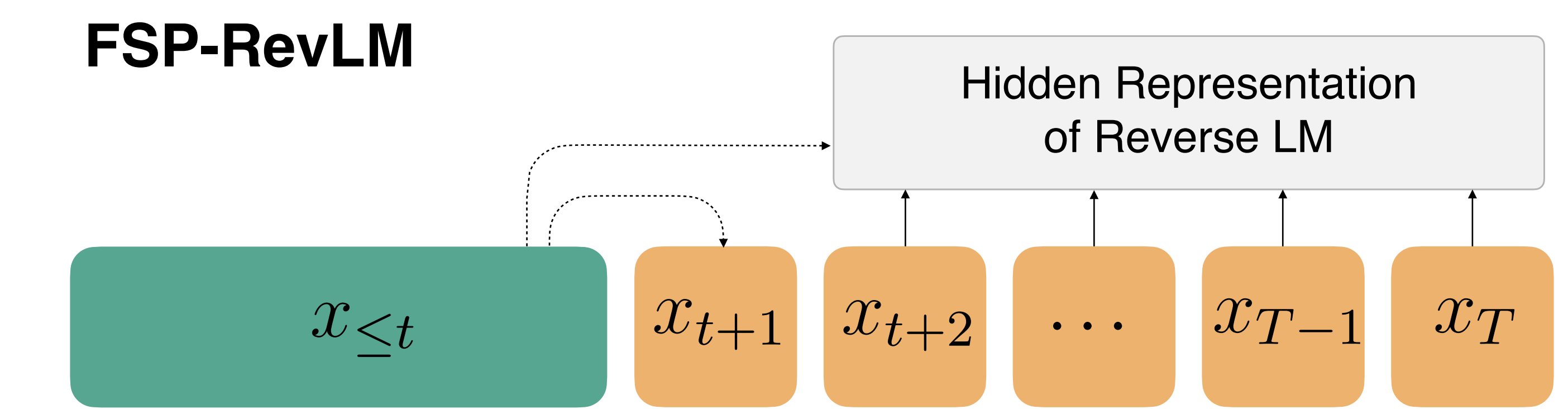
MTP



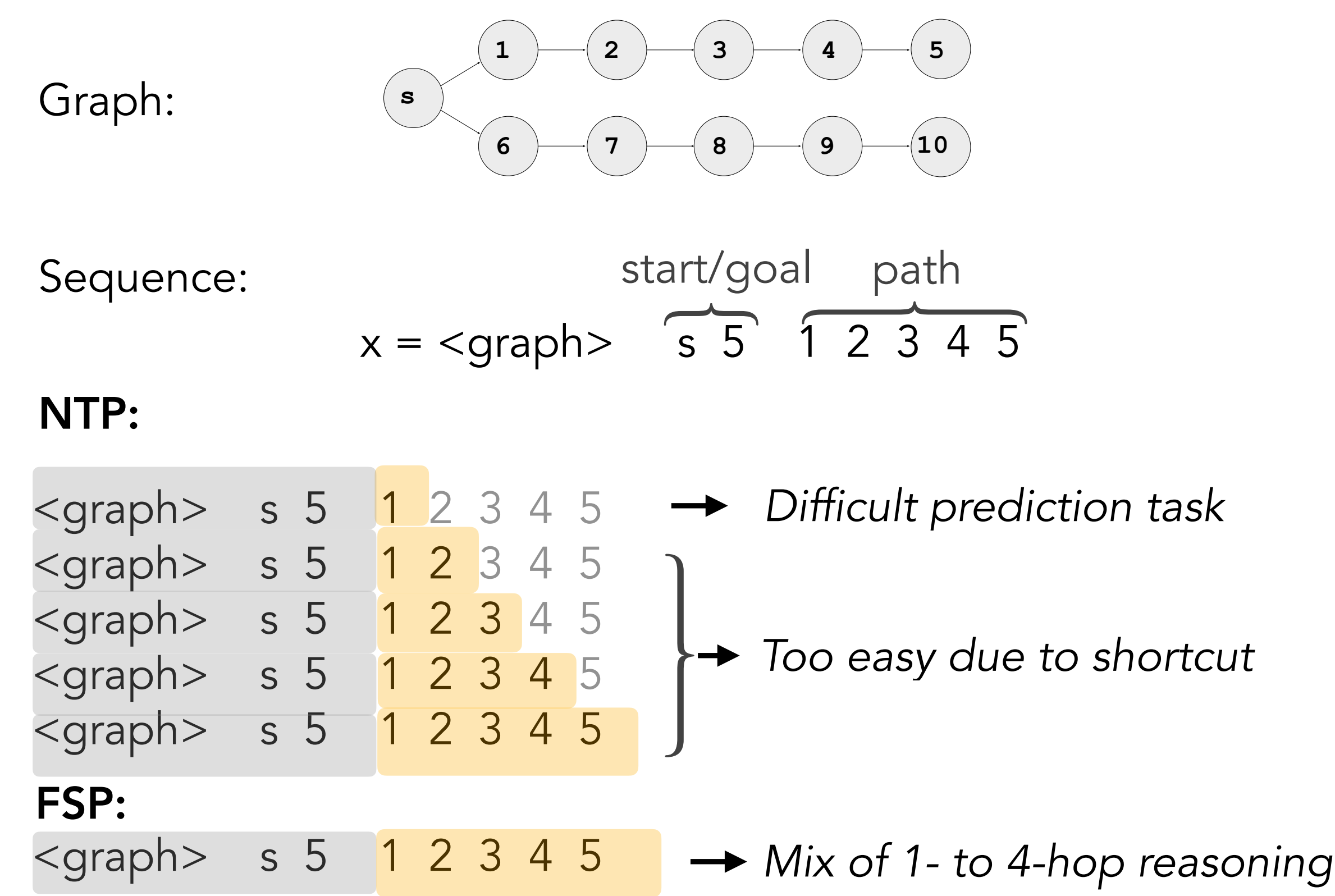
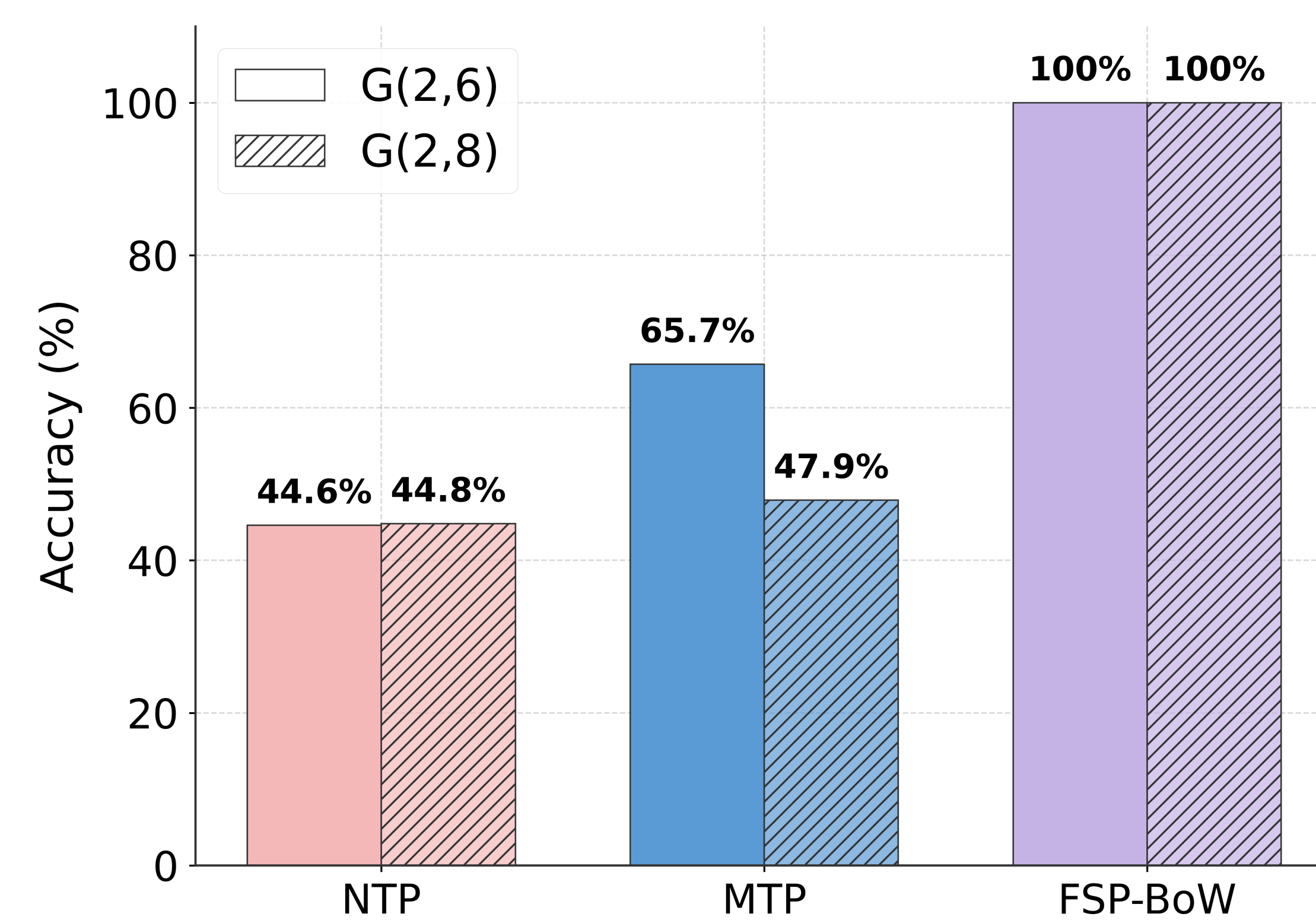
FSP-BoW



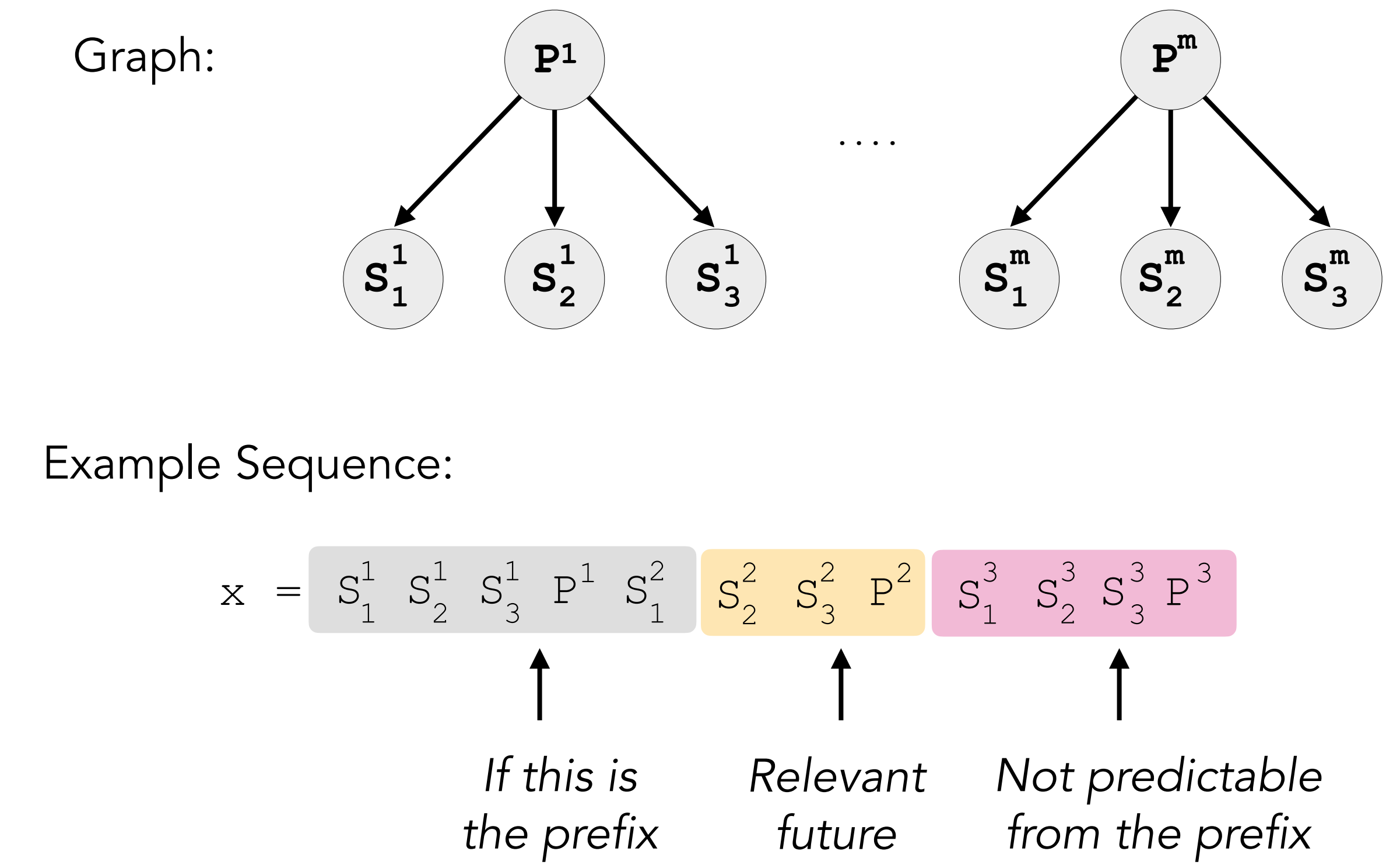
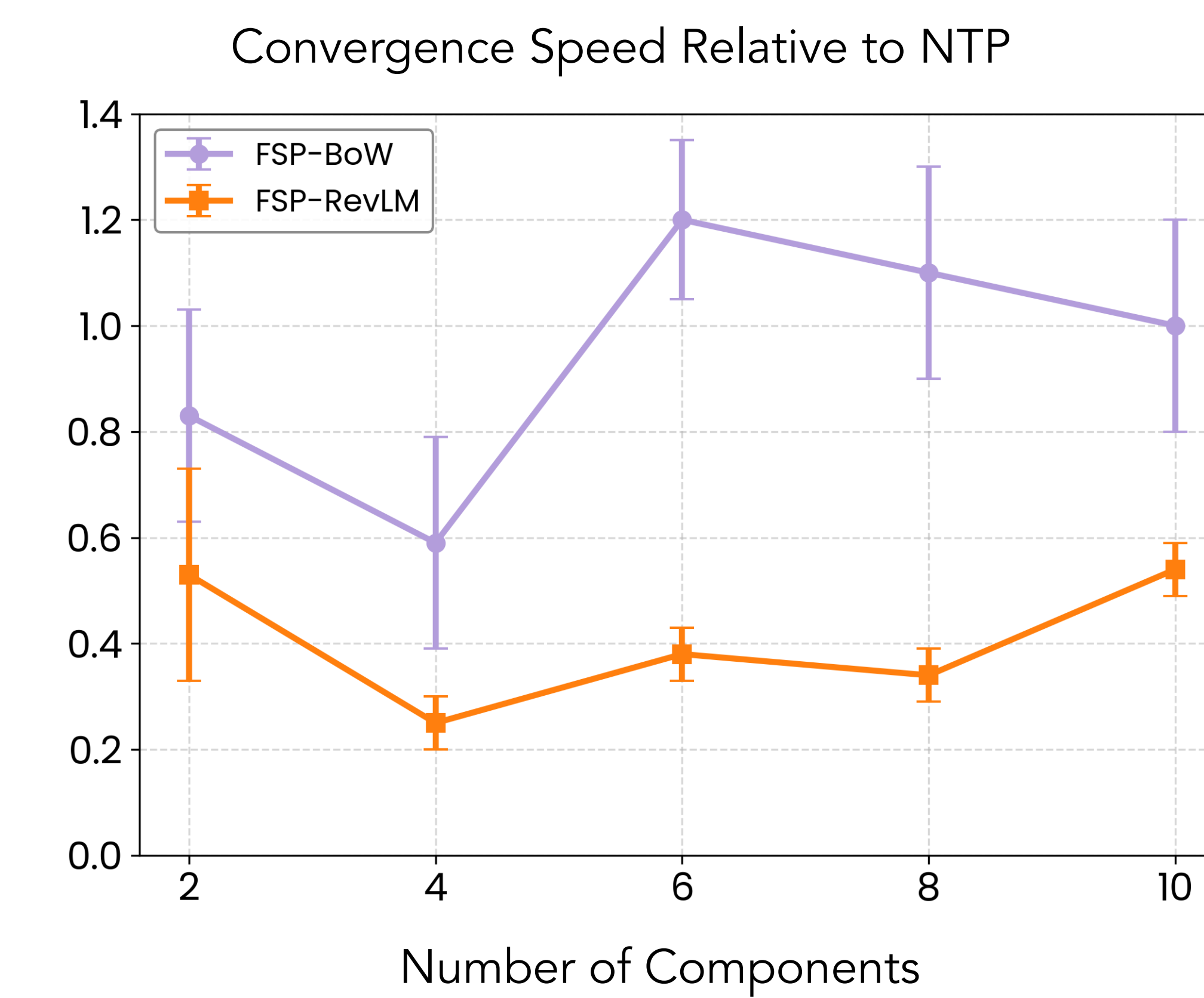
FSP-RevLM



Long Future Summary is Important



Adaptive Future Summary is Important



Pretraining Results: 8B Scale

