

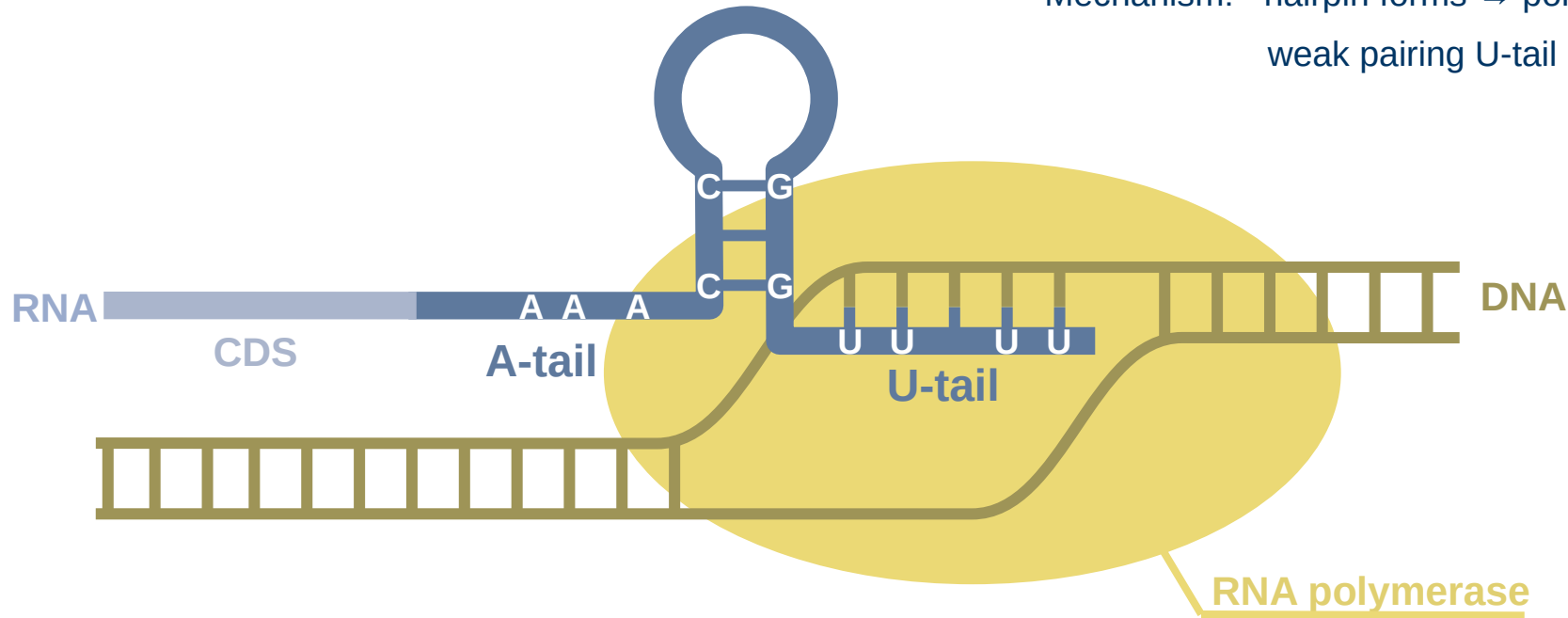
Inverse folding based pre-training for the reliable Identification of intrinsic transcription terminators

Vivian Brandenburg

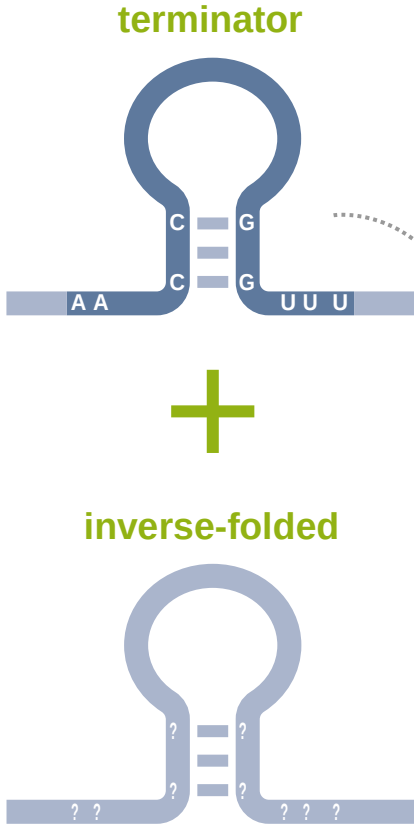
STUDATA 2021

Intrinsic Transcription Terminators

- RNA element on end of nascent transcripts
- Motif: sequence (A-/U-tail) + structure (hairpin)
- Mechanism: hairpin forms → polymerase pauses
weak pairing U-tail → RNA escapes



Data: inverse-folding



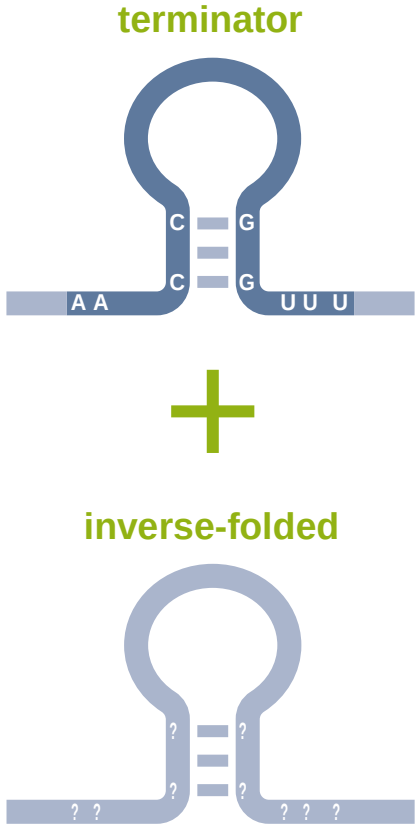
- 1175 terminators from *B. subtilis*¹ and *E. coli*²
- Used for training of Neural Network
- Limited data availability

- Additional data with same structure but different sequence
- Gained with `inverse_fold`³

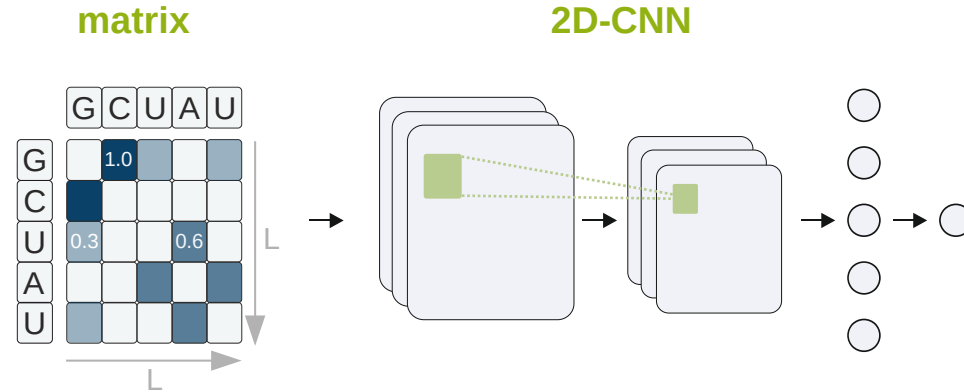
terminator structure . . (((. . .))) . .
random sequence n n n n n n n n n n n n
inverse-fold sequence n n **U C C** n n n **G G A** n n

- **Pre-train** with inverse folded data first, then with terminators

Data



Topology

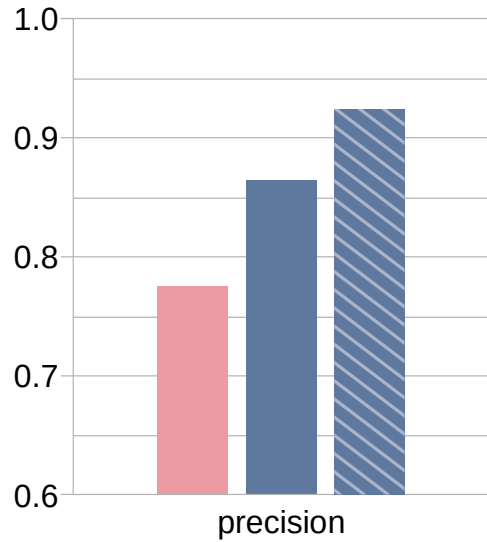


models

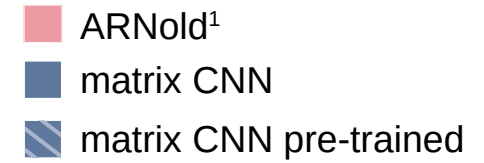
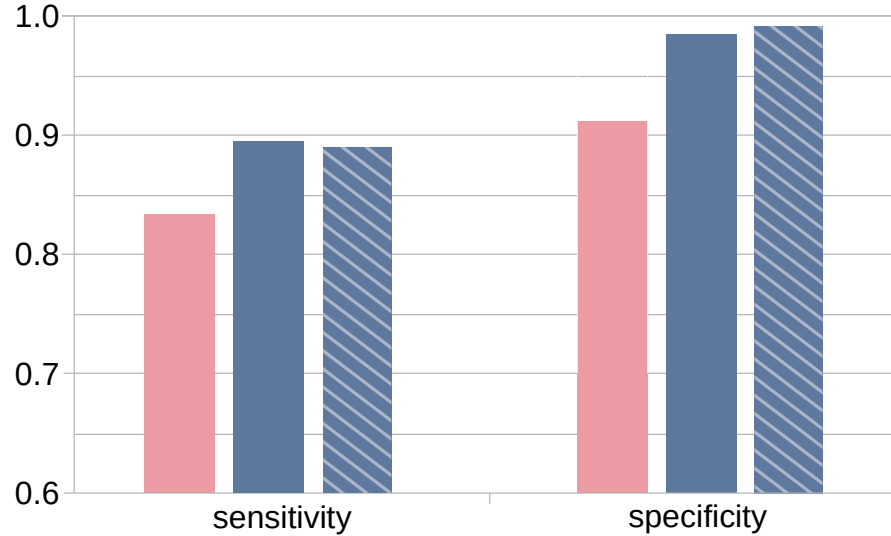
- matrix CNN
- matrix CNN pre-trained

Performance

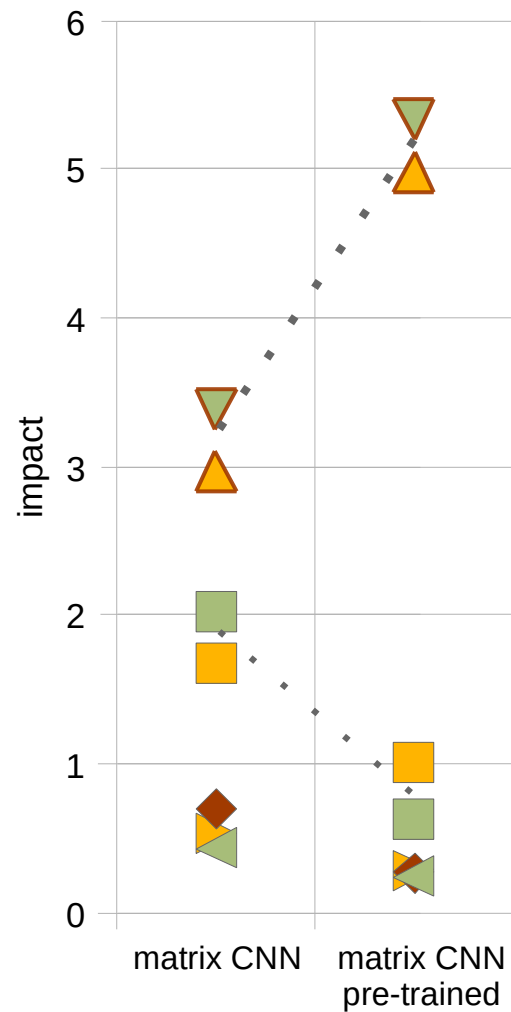
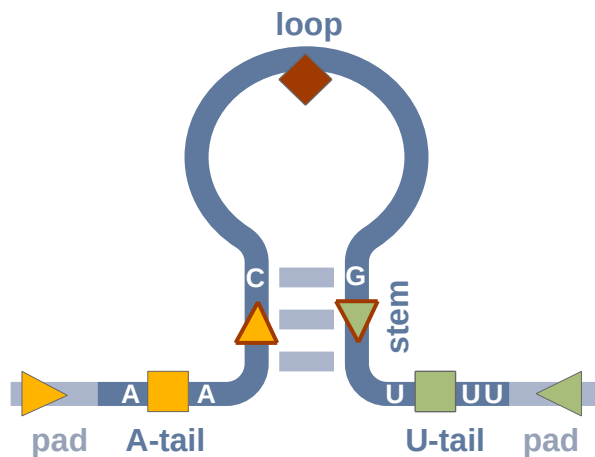
E. coli whole genome



test dataset



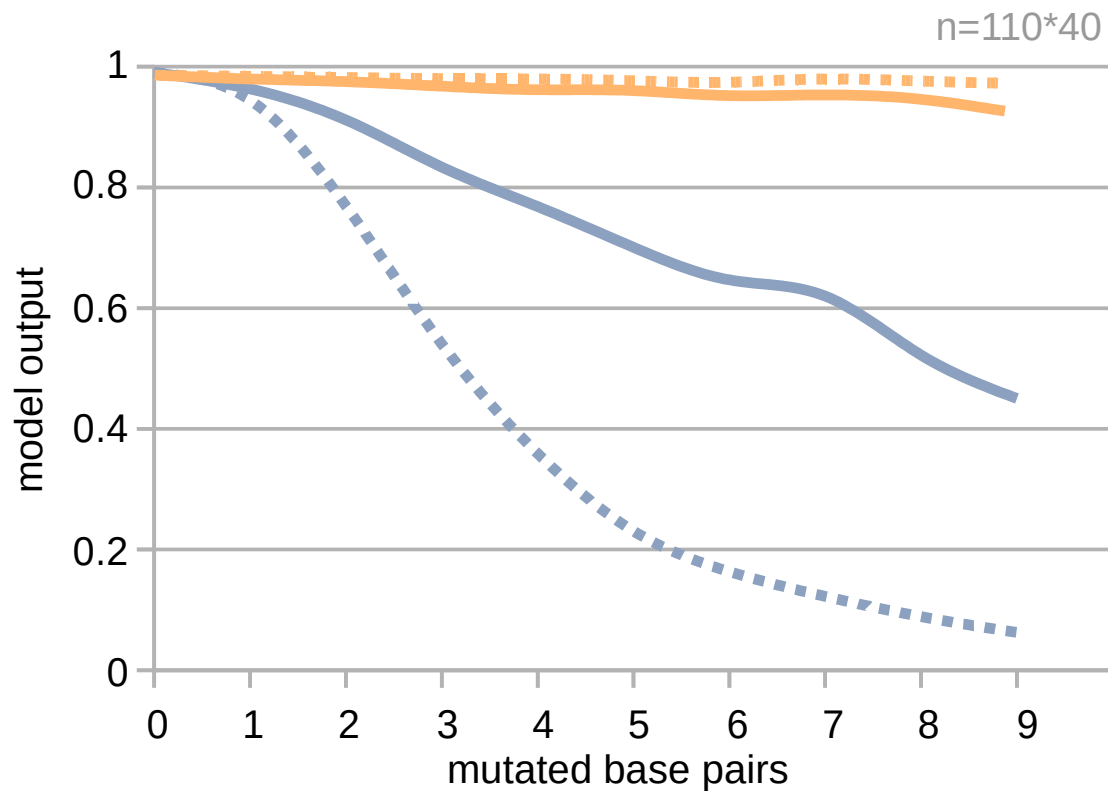
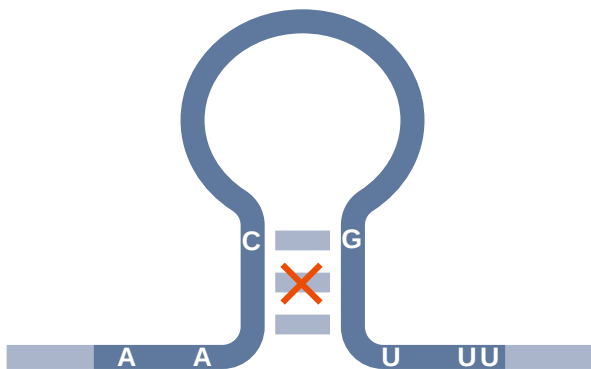
Impact: Sections



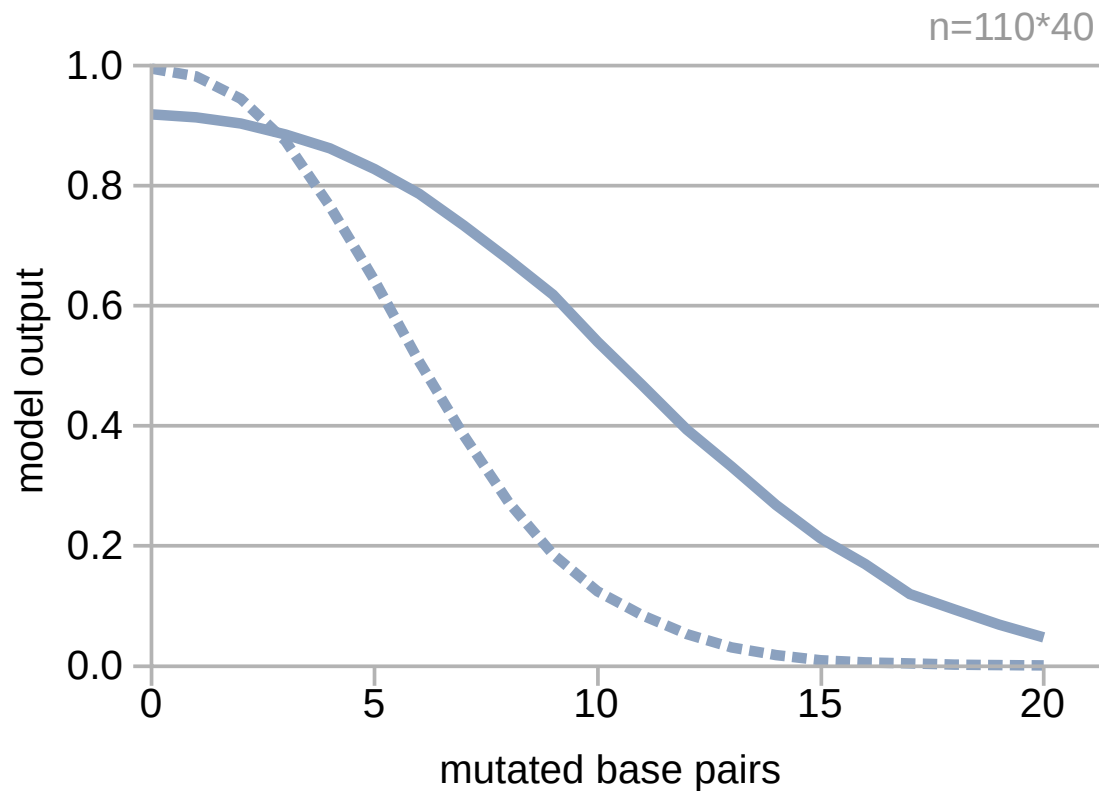
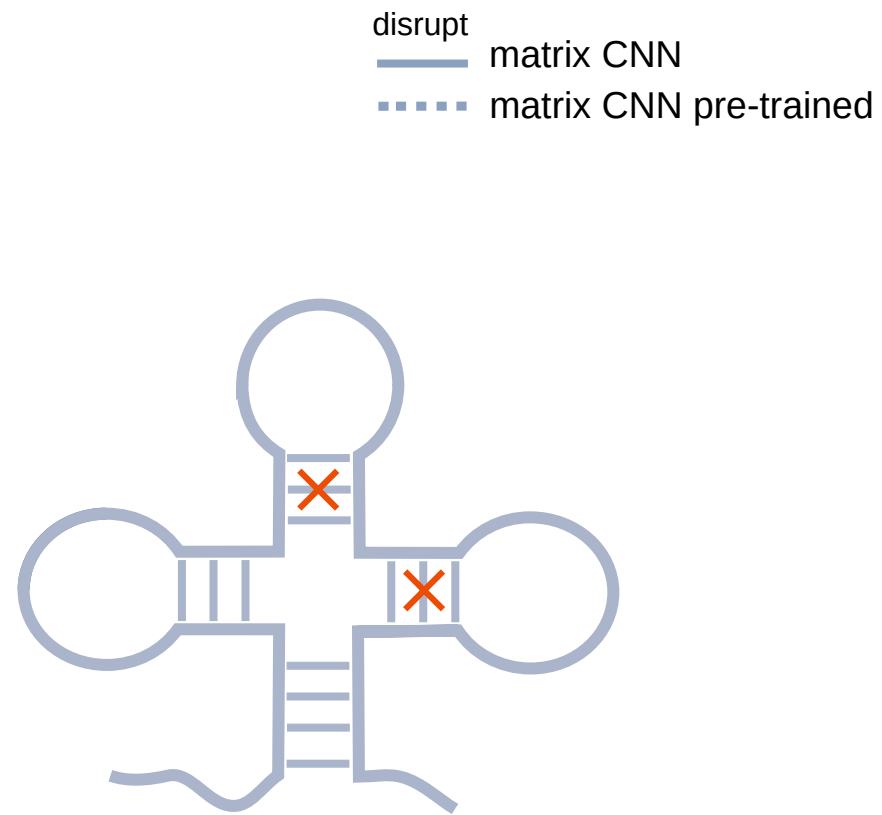
Impact: Base Pairs

preserve disrupt

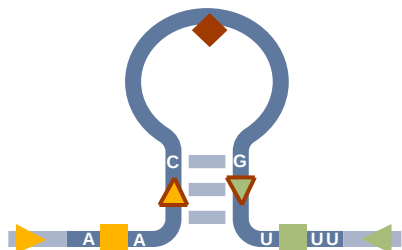
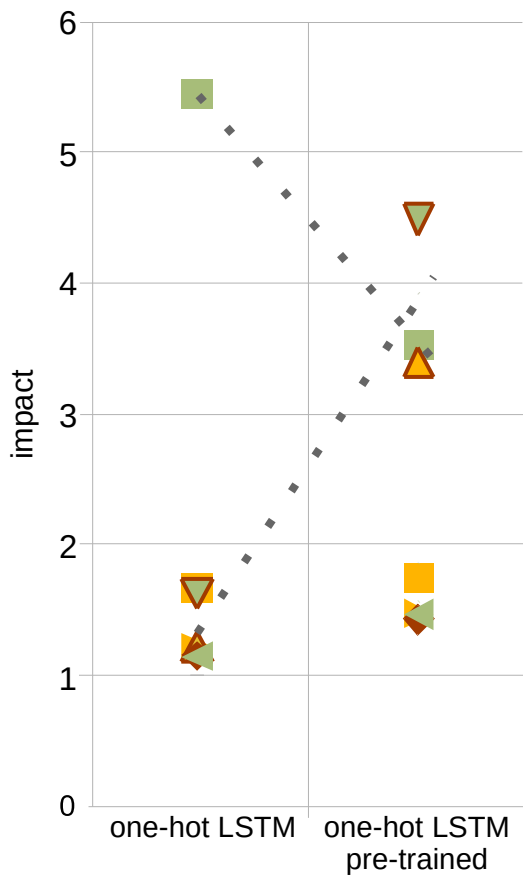
— matrix CNN
- - - matrix CNN pre-trained



Application to other RNAs

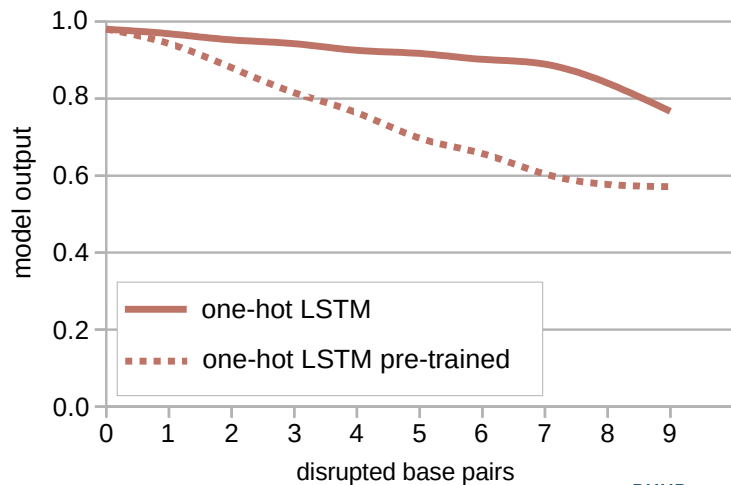
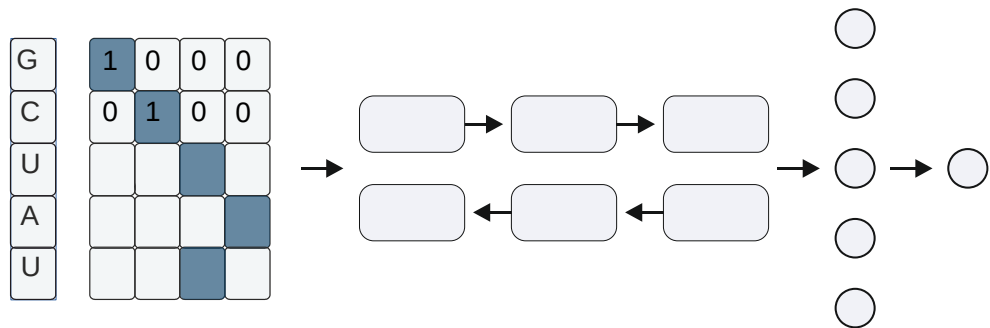


Application to other Topologies



one-hot

LSTM



Summary:

- Improves detection of Intrinsic Transcription Terminators
- Enhances learning of RNA structures
- Is applicable to other RNA families & different deep learning topologies
- Can be used to (partially) overcome limited data availability

