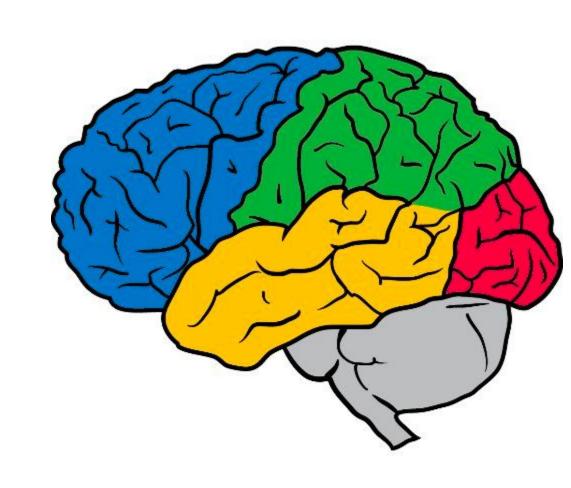


Monotonic Chunkwise Attention

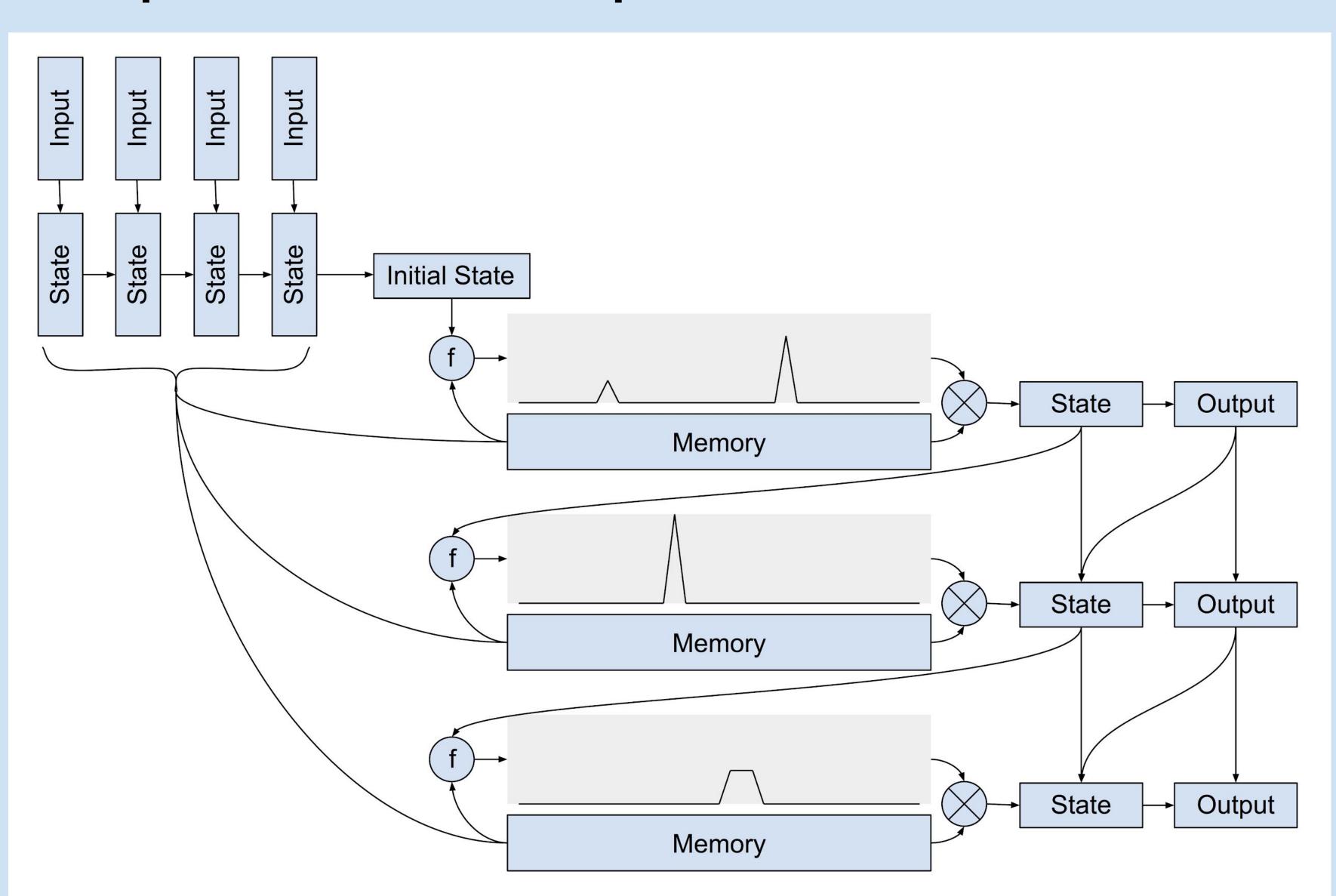


Chung-Cheng Chiu* and Colin Raffel* (equal contribution)

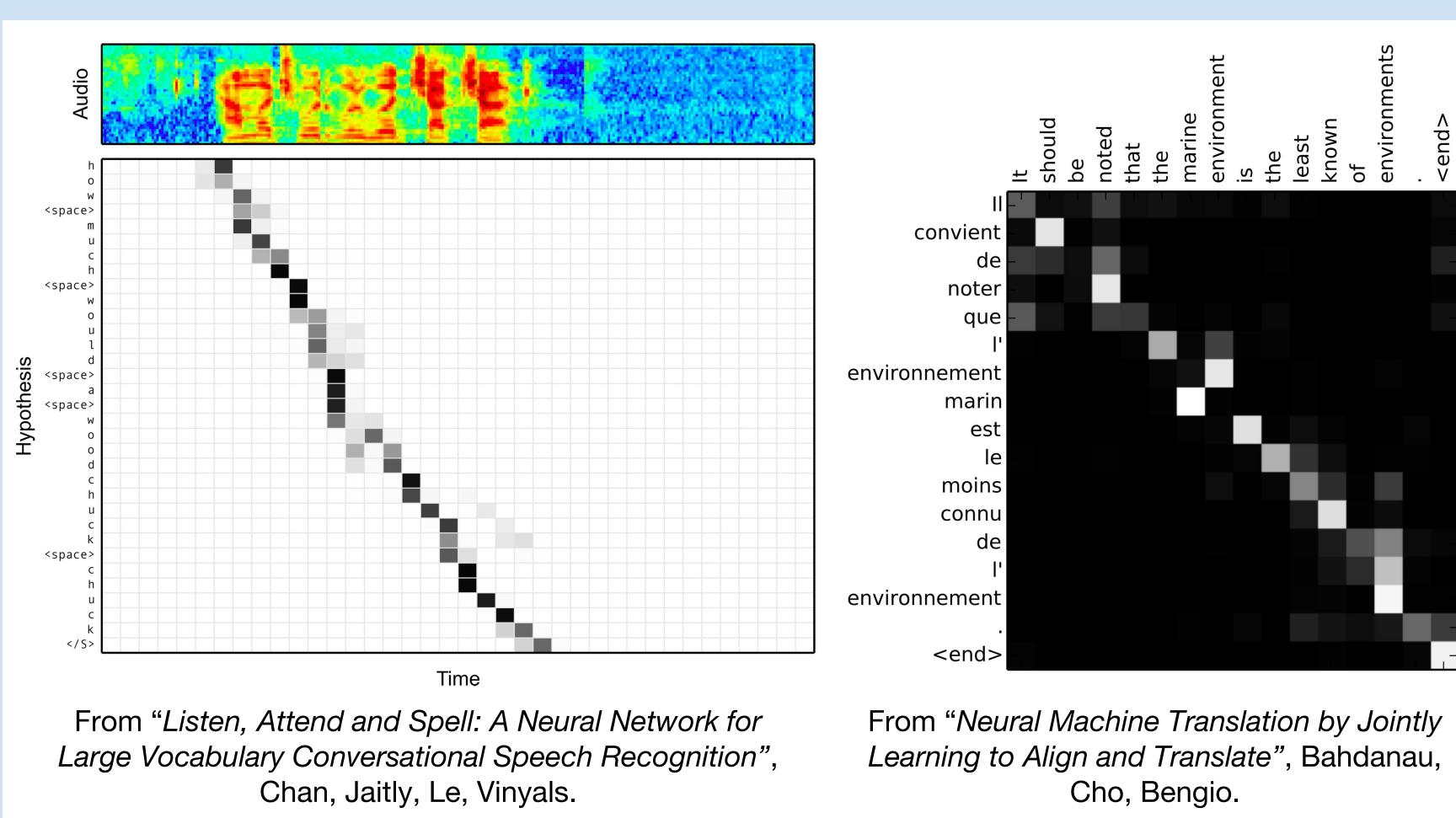
Abstract

Sequence-to-sequence models with soft attention have been successfully applied to a wide variety of problems, but their decoding process incurs a quadratic time and space cost and is inapplicable to real-time sequence transduction. To address these issues, we propose Monotonic Chunkwise Attention (MoChA), which adaptively splits the input sequence into small chunks over which soft attention is computed. We show that models utilizing MoChA can be trained efficiently with standard backpropagation while allowing online and linear-time decoding at test time. When applied to online speech recognition, we obtain state-of-the-art results and match the performance of a model using an offline soft attention mechanism. In document summarization experiments where we do not expect monotonic alignments, we show significantly improved performance compared to a baseline monotonic attention-based model.

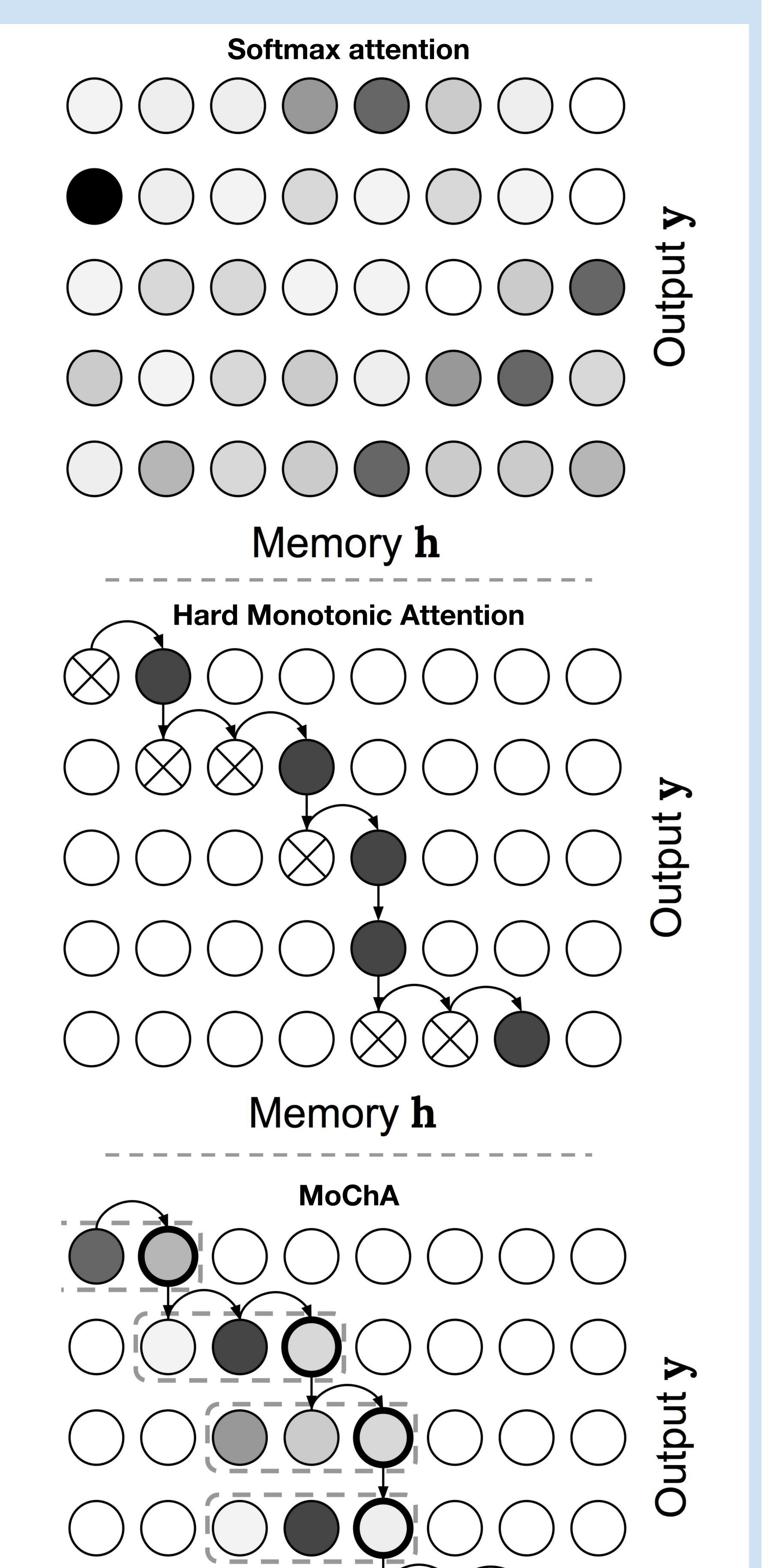
Sequence-to-sequence framework



Attention is often roughly monotonic

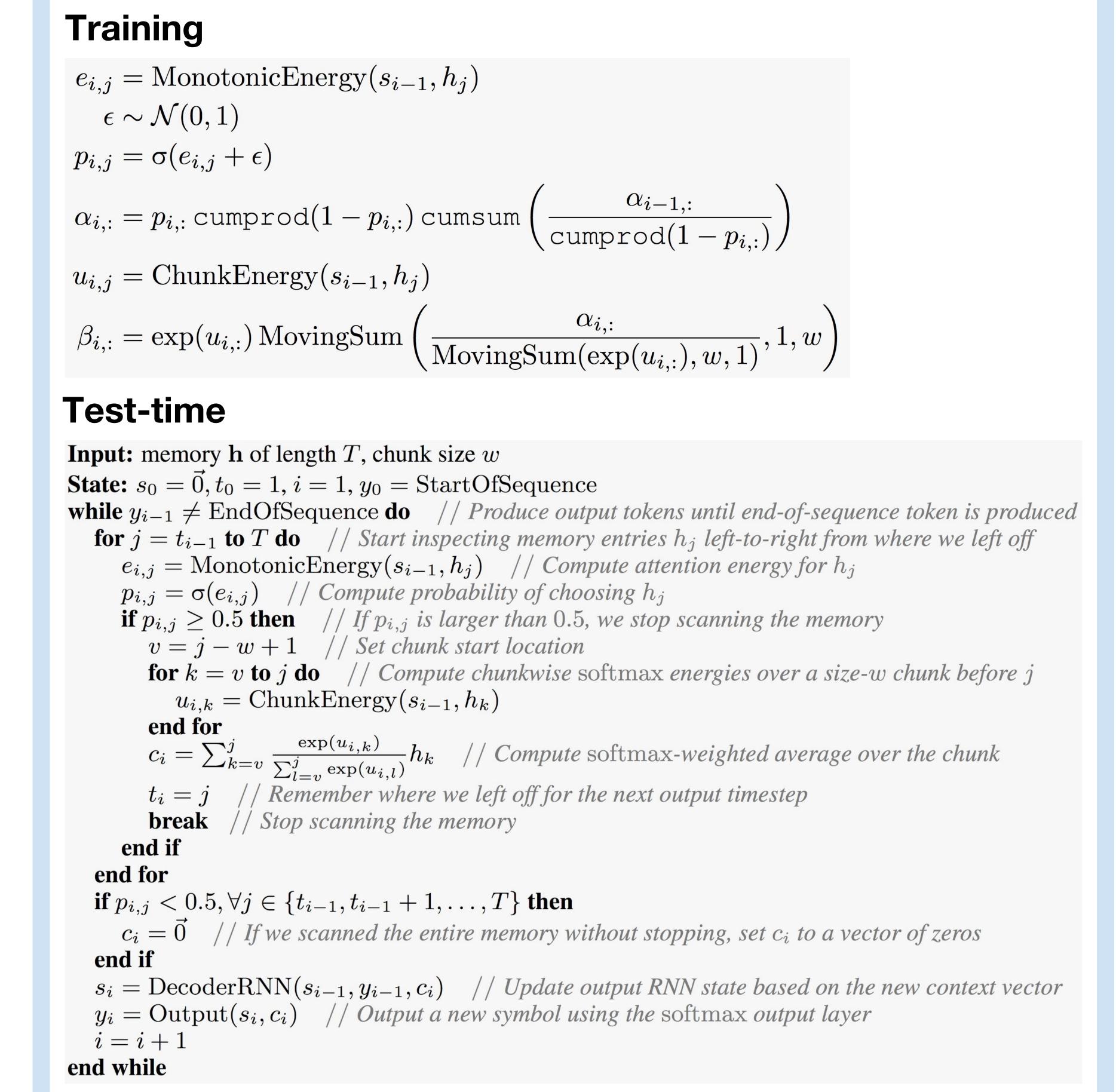


Attention mechanisms



Memory h

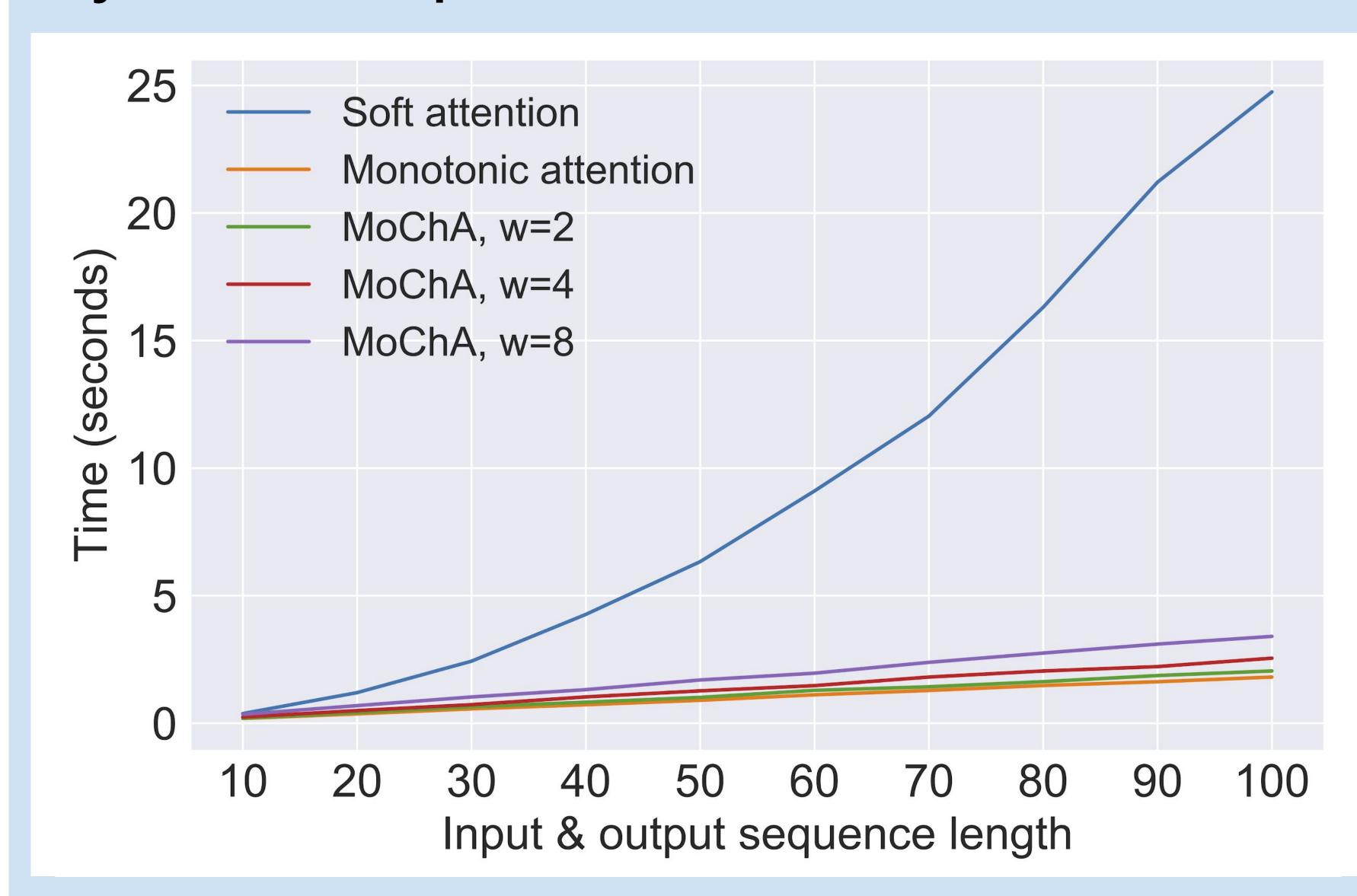
Decoding with MoChA



Results

Prior Result			WER	Mechanism	R-1	R-2
(Luo et al., 2016) (Reinforcement Learning) (Wang et al., 2016) (CTC)			33.4% 27.0% 22.7% 17.4%	Soft Attention (offline) Hard Monotonic Attention MoChA, $w=8$	39.11 31.14 35.46	15.76 11.16 13.55
Attention Mechanism Best WER Average WER			(left) Word Error rate for online speech recognition on WSJ. (above) ROUGE-1			
Soft Attention (offline) MoChA, $w=2$	14.2% 13.9%	$14.6 \pm 0.3\%$ $15.0 \pm 0.6\%$		and 2 for document summarization on CNN/Daily Mail.		

Synthetic speed benchmark



Speech Attention Plots

