

## Abstract

We demonstrate a generic speech engine that handles a broad range of scenarios without needing to resort to domain-specific optimizations.

We perform a focused search through model architectures finding deep recurrent nets with multiple layers of 2D convolution and layer-to-layer batch normalization to perform best.

In many cases, our system is competitive with the transcription performance of human workers when benchmarked on standard datasets.

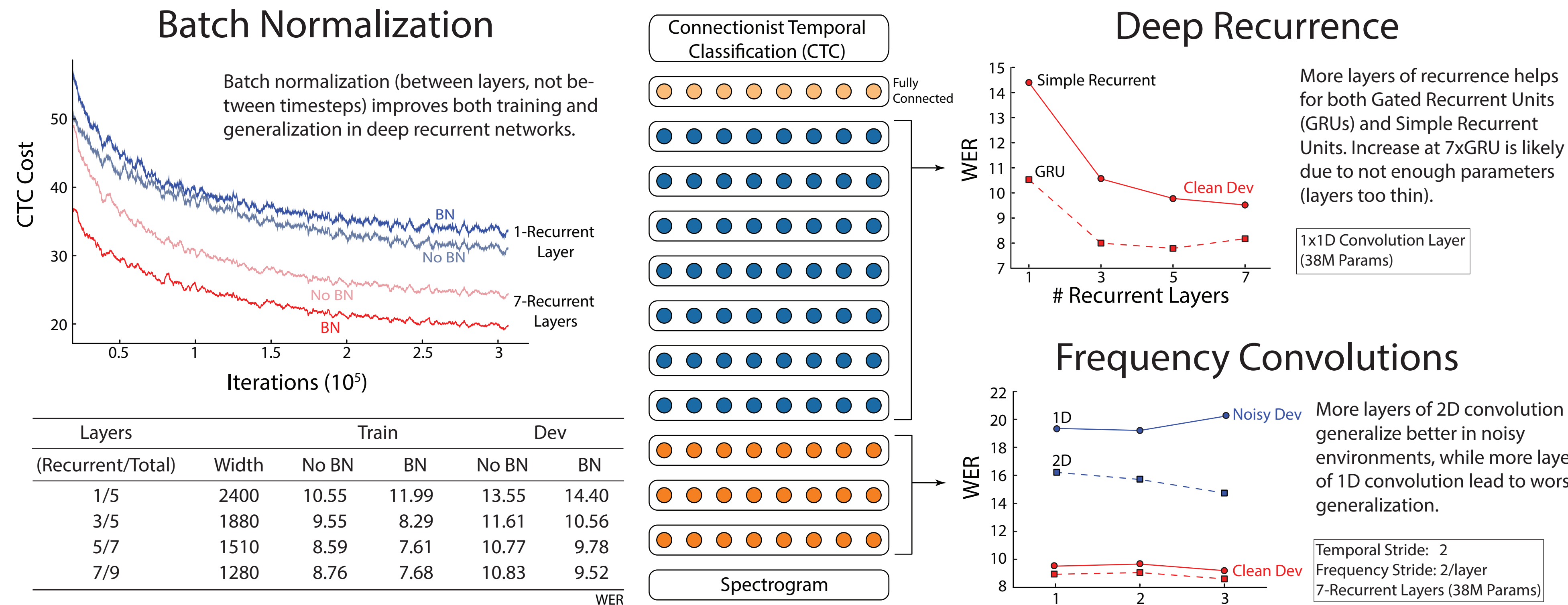
This approach also adapts easily to multiple languages, which we illustrate by applying essentially the same system to both Mandarin and English speech.

To train on 12,000 hours of speech, we perform many systems optimizations. Placing all computation on distributed GPUs, we can perform synchronous SGD at 3 TFLOP/s per GPU for up to a 128 GPUs (weak scaling).

For deployment, we develop a batching scheduler to improve computational efficiency while minimizing latency. We also create specialized matrix multiplication kernels to perform well at small batch sizes.

Combined with forward only variants of our research models, we achieve a low-latency production system with little loss in recognition accuracy.

## Architecture Search / Data Collection



### SortaGrad

	Baseline	Batch Norm
Not Sorted	11.96	9.78
Sorted	10.83	9.52

WER

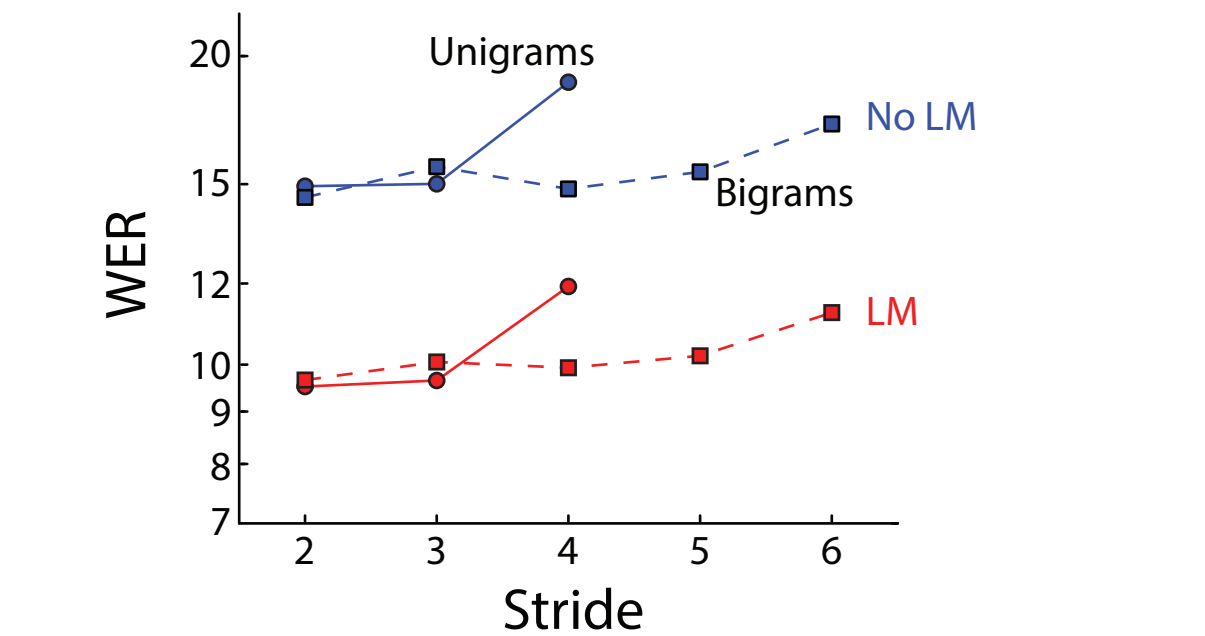
"A Sort of Stochastic Gradient Descent" Curriculum learning with utterance length as a proxy for difficulty. Sorting the first epoch by length improves convergence stability and final values.

### Diverse Datasets

Dataset	Speech Type	Hours
WSJ	read	80
Switchboard	conversational	300
Fisher	conversational	2000
LibriSpeech	read	960
Baidu	read	5000
Baidu	mixed	3600
Total		11940

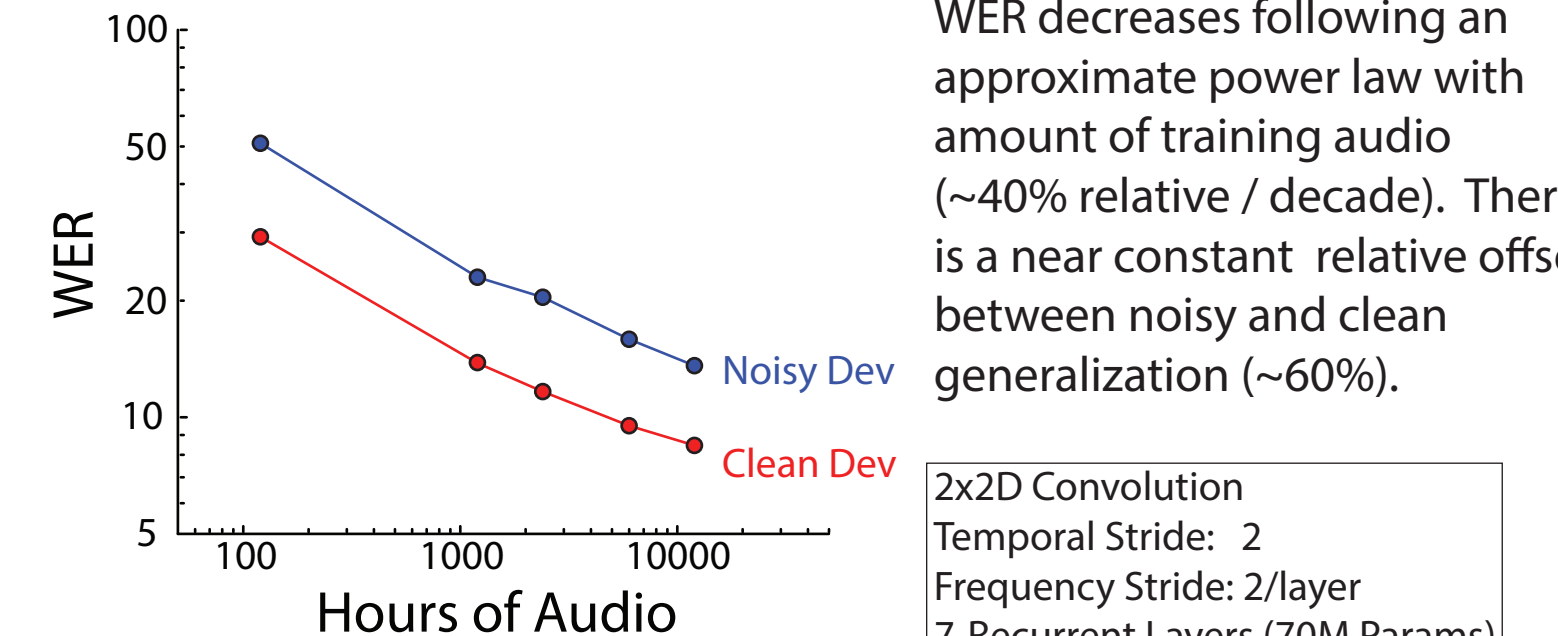
Mix of conversational, read, and spontaneous. Average utterance length ~6 seconds.

### Bigrams

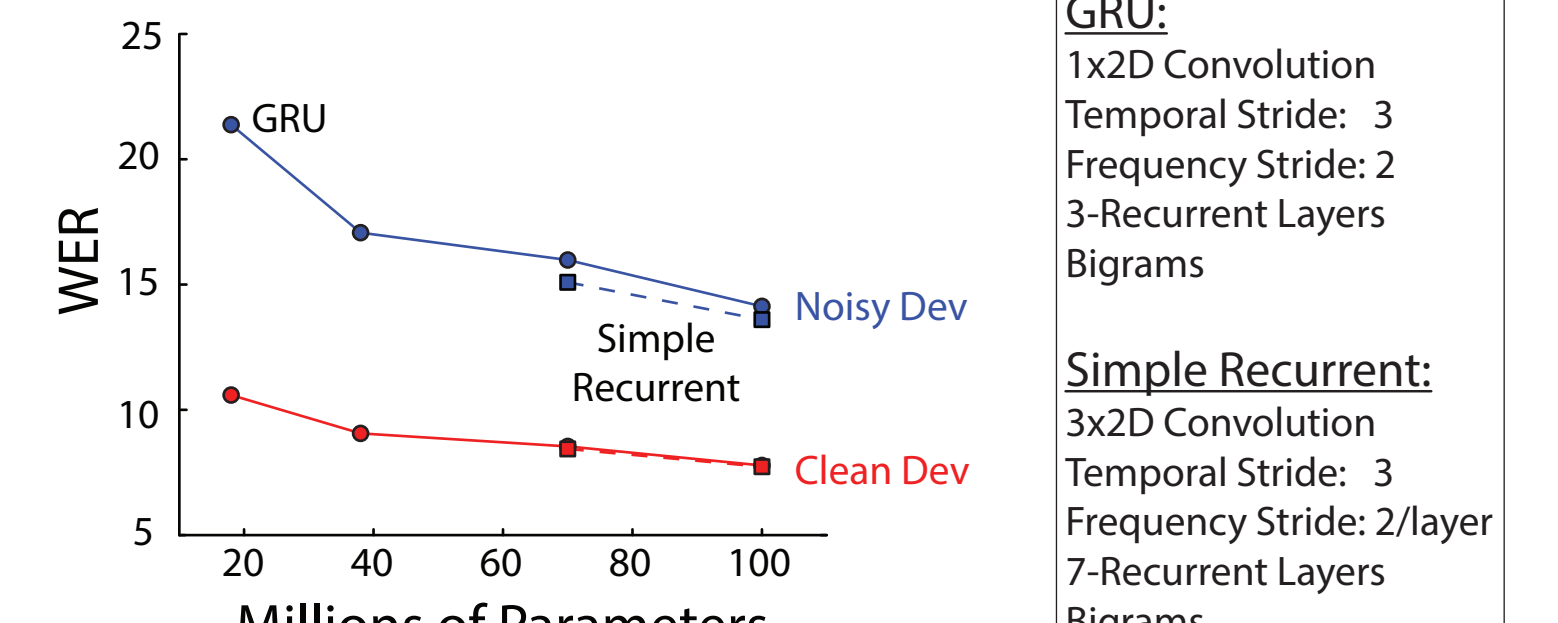


	Unigram	Mandarin	Bigram
Chars / Sec	14.1	3.2	2.4
Bits / Char	4.1	12.6	9.6
Bits / Sec	58.1	40.7	23.3

### Data Scaling



### Model Size



## Single Model Approaching Human Performance on Many Test Sets

Read Speech				Accented Speech				Noisy Speech			
Test set	DS1	DS2	Human	Test set	DS1	DS2	Human	Test set	DS1	DS2	Human
WSJ eval'92	4.94	3.60	5.03	VoxForge American-Canadian	15.01	7.55	4.85	CHiME eval clean	6.30	3.34	3.46
WSJ eval'93	6.94	4.98	8.08	VoxForge Commonwealth	28.46	13.56	8.15	CHiME eval real	67.94	21.79	11.84
LibriSpeech test-clean	7.89	5.33	5.83	VoxForge European	31.20	17.55	12.76	CHiME eval sim	80.27	45.05	31.33
LibriSpeech test-other	21.74	13.25	12.69	VoxForge Indian	45.35	22.44	22.15				

DS 1	DS 2	Human
Human et al., Deep speech: Scaling up end-to-end speech recognition, 2014. <a href="http://arxiv.org/abs/1412.5567">http://arxiv.org/abs/1412.5567</a>	1x1D Convolution Temporal Stride: 2 Unigrams / No BN 1-Recurrent Layer (38M Params)	3x2D Convolution (Bigrams / BN) Temporal Stride: 3 Frequency Stride: 2/layer 7-Recurrent Layers (100M Params)

## Architectural Improvements Generalize to Mandarin ASR

Dataset	Type	Hours	Language	Architecture	Dev No LM	Dev LM	Architecture	Dev	Test
Unaccented Mandarin	spontaneous	5000	English	5-layers, 1 recurrent	27.79	14.39	5-layers, 1 recurrent	7.13	15.41
Accented Mandarin	spontaneous	3000	English	9-layers, 7 recurrent	14.93	9.52	5-layers, 3 recurrent	6.49	11.85
Assorted Mandarin	mixed	1400	Mandarin	5-layers, 1 recurrent	9.80	7.13	5-layers, 3 recurrent + BN	6.22	9.39
Total		9400	Mandarin	9-layers, 7 recurrent	7.55	5.81	9-layers, 7 recurrent + BN + 2D conv	5.81	7.93

Mix of conversational, read, and spontaneous. Average utterance length ~3 seconds.

WER English / CER Mandarin

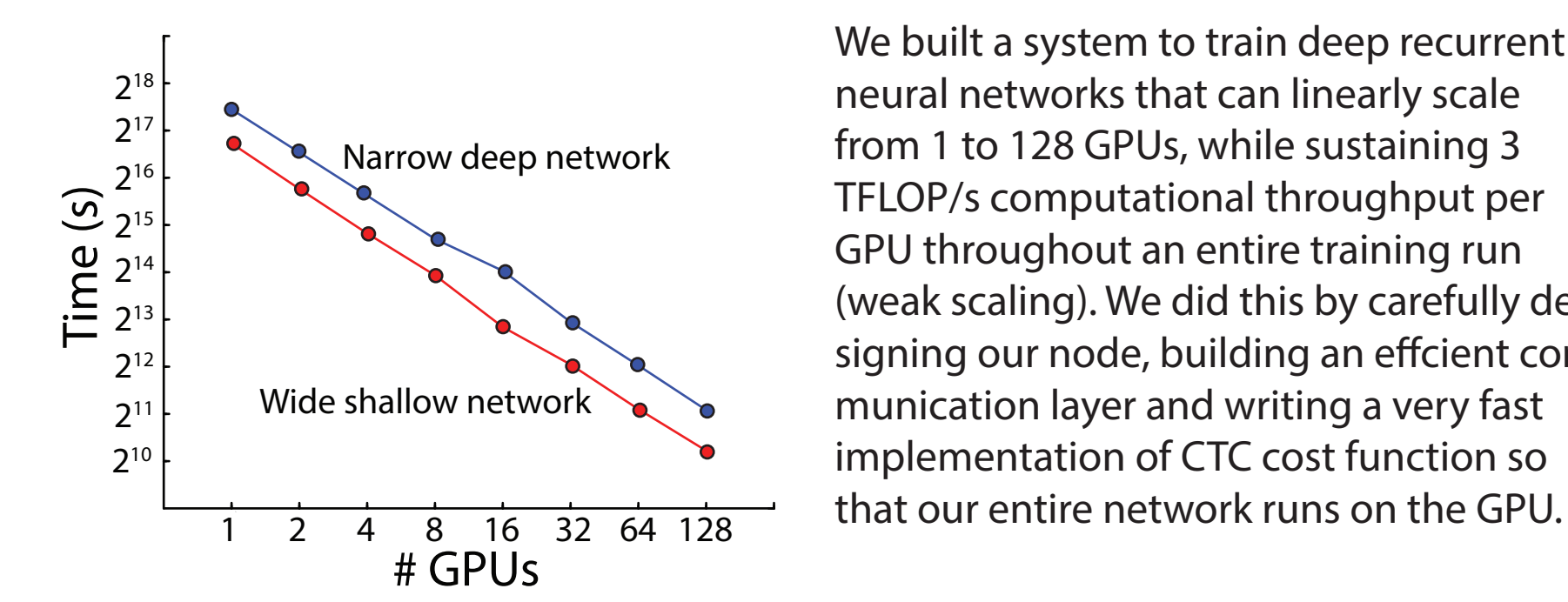
3x2D Convolution  
Temporal Stride: 4  
Frequency Stride: 2/layer  
70M Params

Internal Mandarin voice search test set.

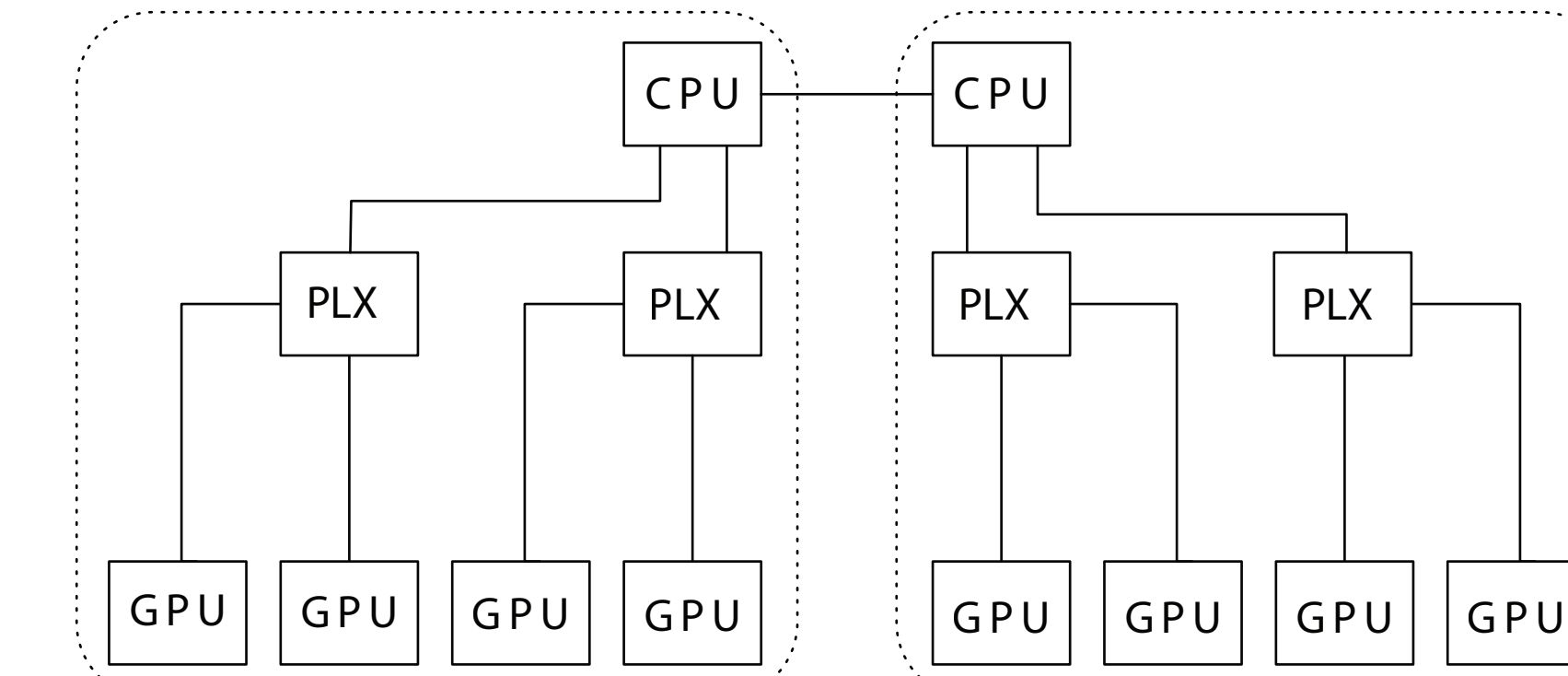
CER

## Systems Optimizations / Deployment

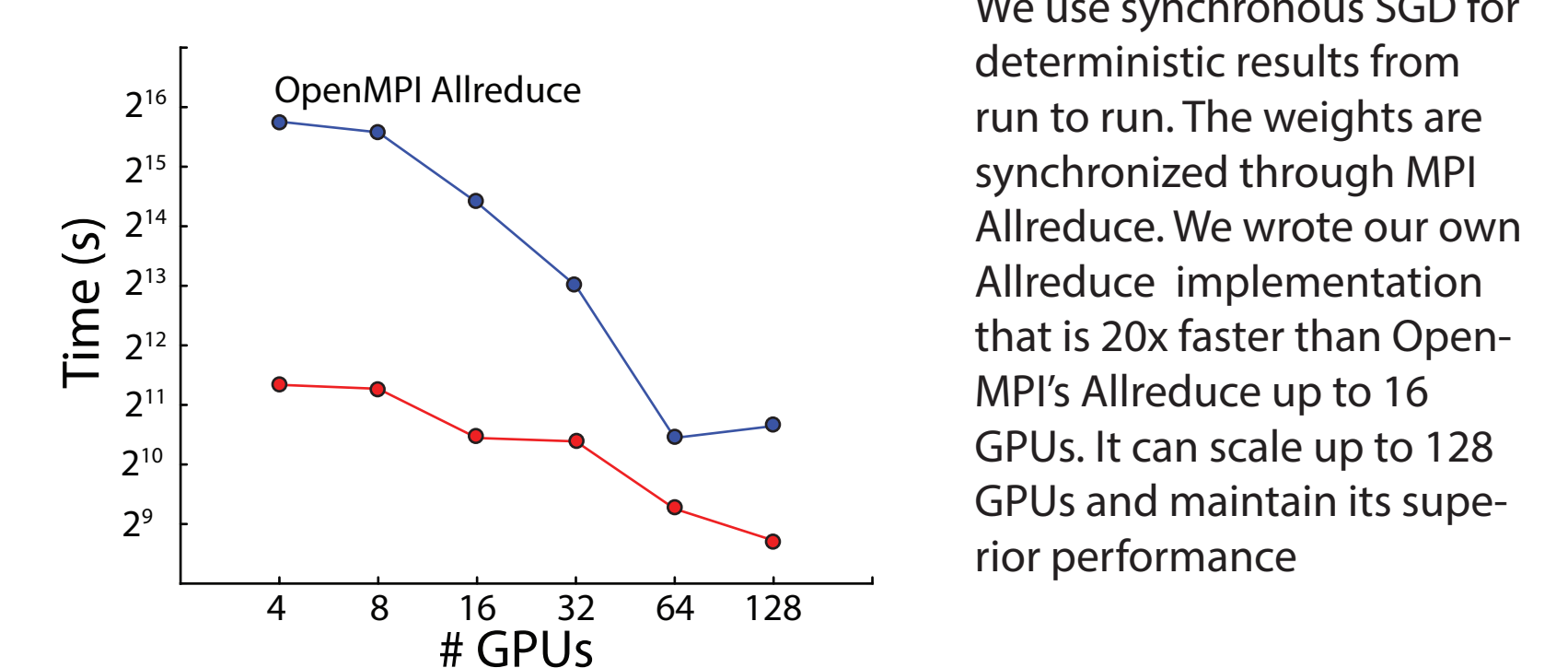
### Scaling Neural Network Training



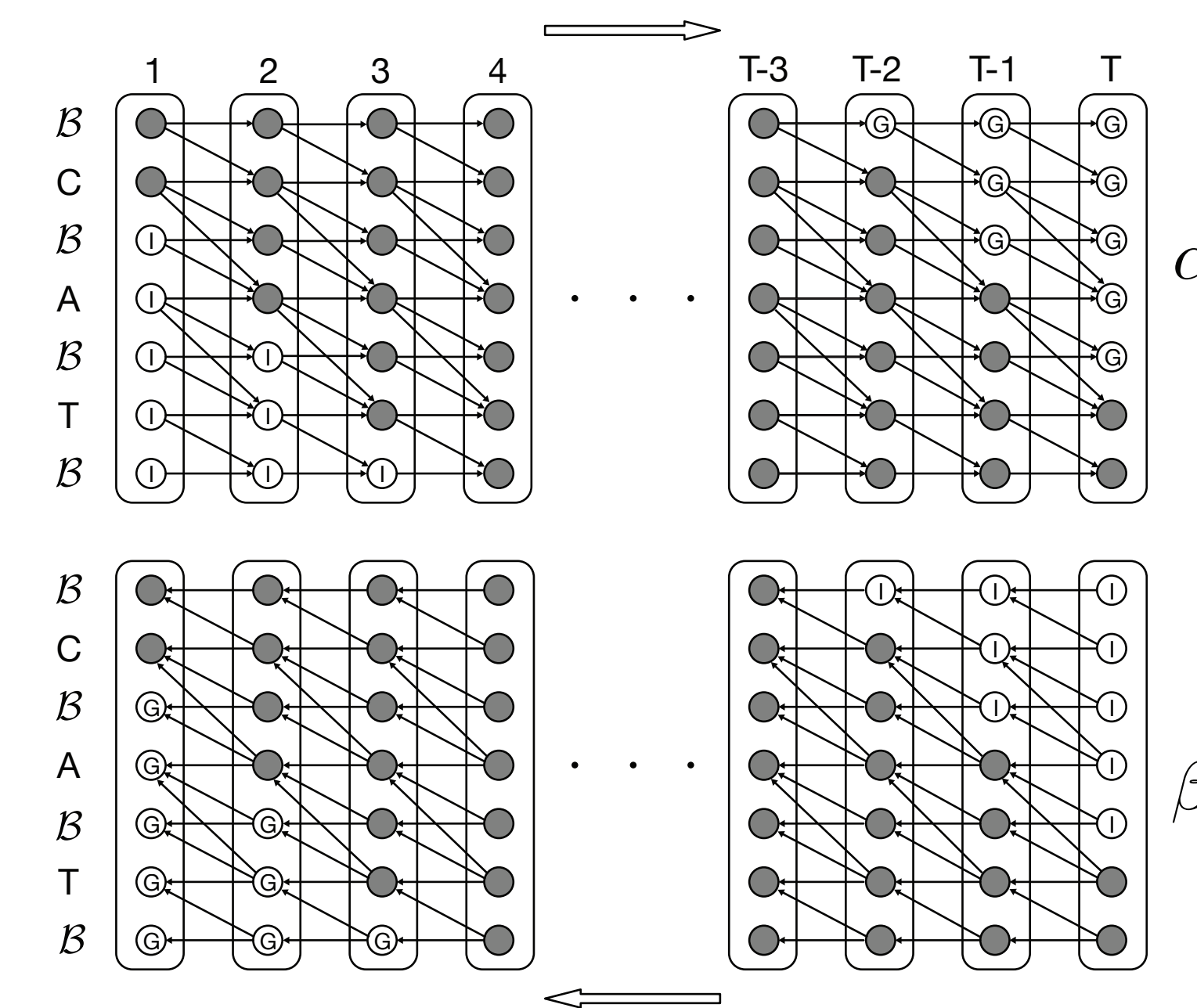
### Node Architecture



### Fast MPI AllReduce



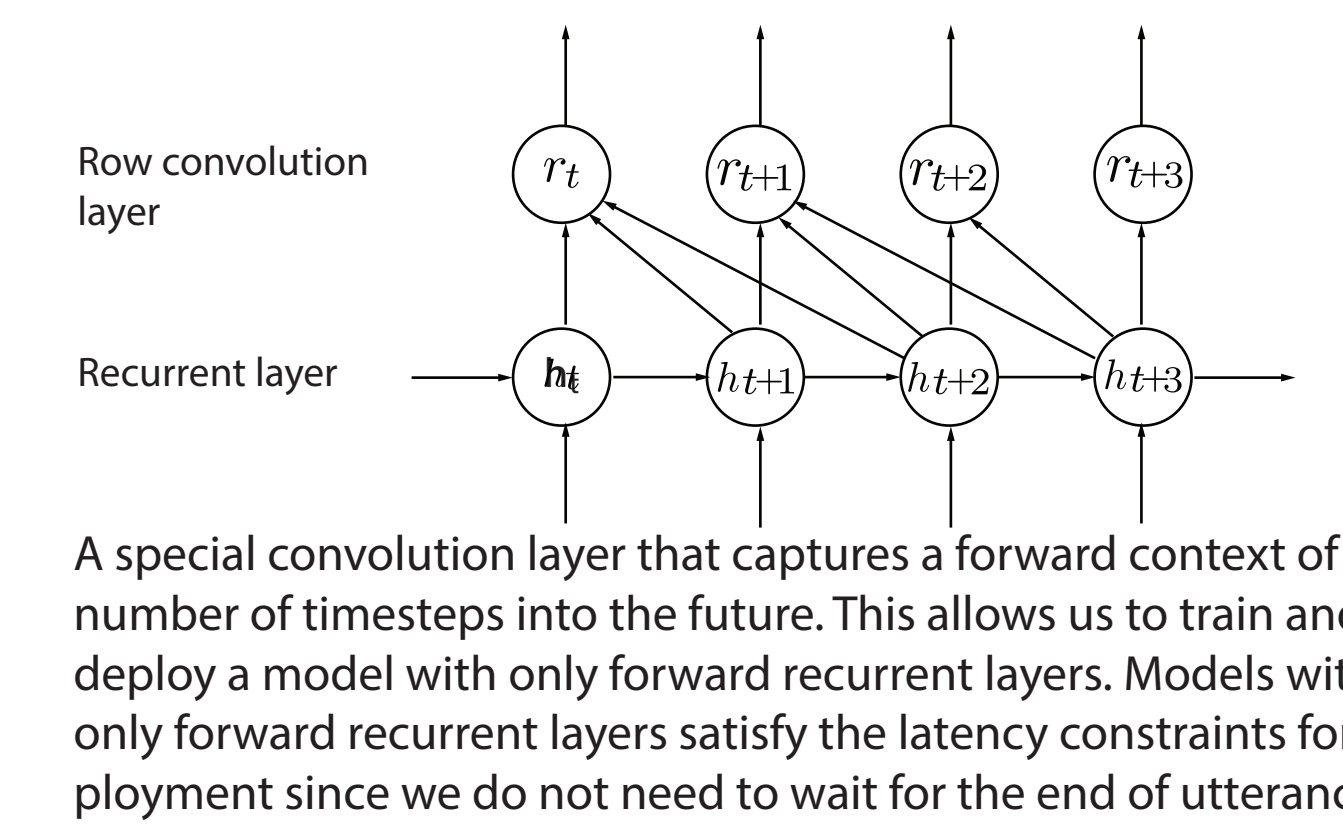
### Computing CTC on GPU



Language	CPU CTC Time	GPU CTC Time	Speedup
English	5888.12	203.56	28.9
Mandarin	1688.01	135.05	12.5

An optimized CTC implementation on the GPU is significantly faster than the corresponding CPU implementation. This also eliminates all roundtrip copies between CPU and GPU. All times are in seconds for training one epoch of a 5 layer, 3 bi-directional layer network.

### Row Convolutions

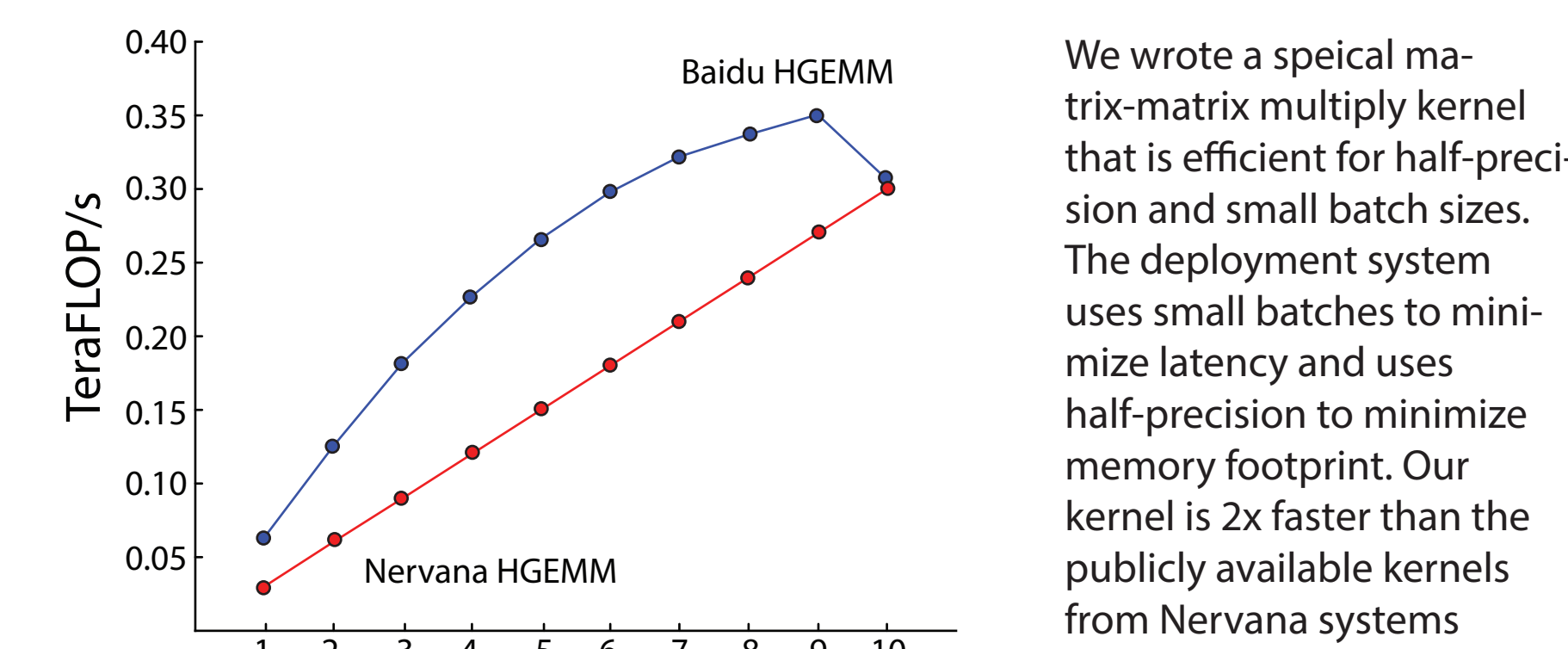


As we see on the right, the production system that only has forward recurrent layers is very close in performance to the research system with bi-directional layers, on two different test sets.

System	Clean	Noisy
Research	6.05	9.75
Production	6.10	10.00

CER

## Deployment Optimized Matrix-Matrix Multiply



### Batch Dispatch

