



# LSTM Deep Dive

Presented By

Jingxiu Liu, Director of Product Marketing, AI and Edge Computing  
2018/10/1

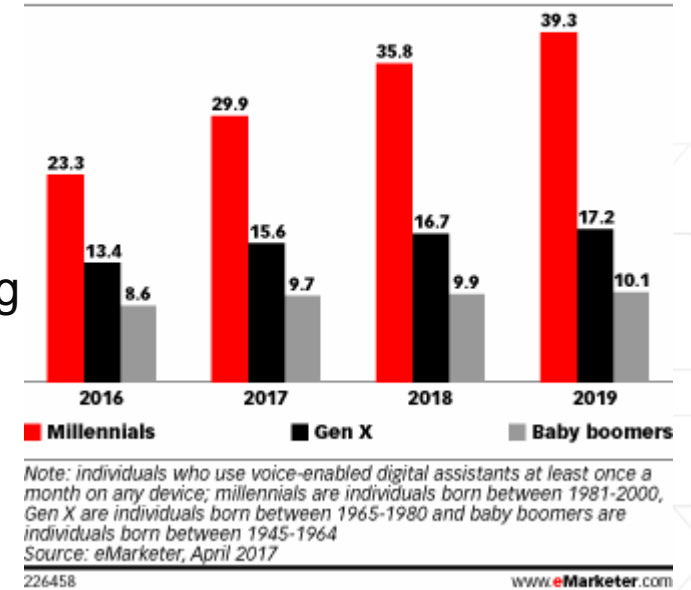


# Speech to Text Application

## How it drive business values?

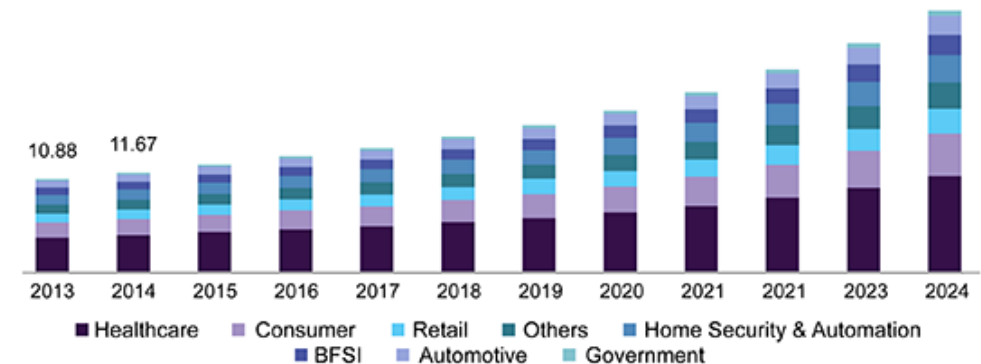
- Usability : the accurate is reported on par with humans
  - the word error rate for **Microsoft's 5.1%**, while **Google 4.9%**.
- virtual assistants with speech recognition capabilities keeps increasing
- Change the way we interact with and build electronics
  - Voicebox : worked on voice recognition for partners including Samsung, AT&T, and Toyota.

**US Voice-Enabled Digital Assistant Users, by Generation, 2016-2019**  
millions



The global voice recognition market size was valued at **USD 55.17 billion in 2016** and is expected to grow at a CAGR of **11.0%** during the forecast period of 2016 to 2024.

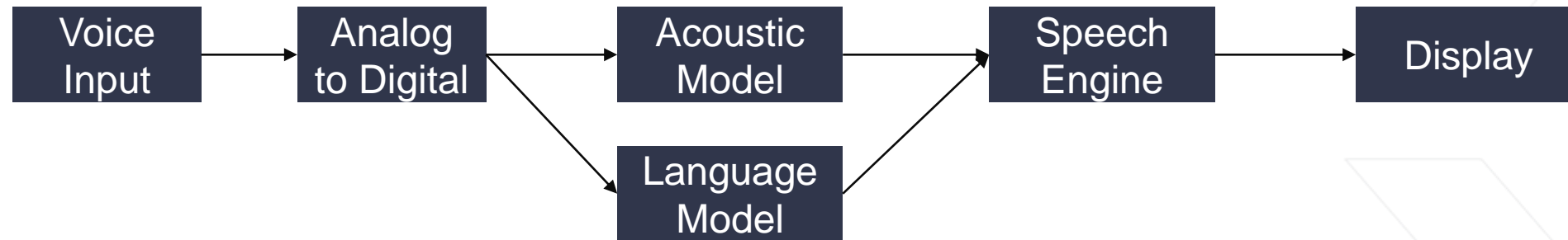
**Europe voice recognition market size, by vertical, 2013 - 2024 (USD Billion)**



# Speech to Text Application

## What is Speech to Text (Speech Recognition) ?

- understand voice by the computer and performing any required task.



## What can it be used ?

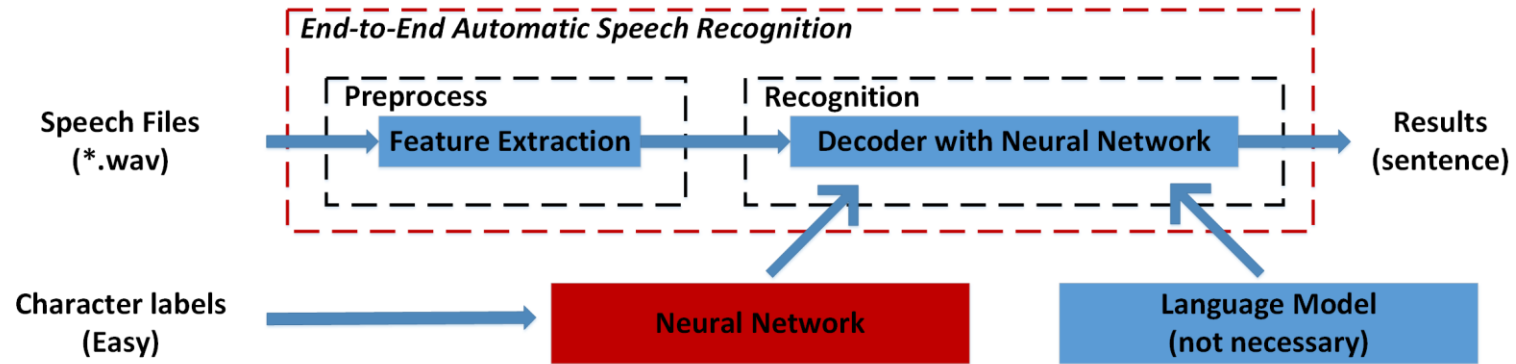
- Dictation
- System control/navigation
- Voice dialing
- Commercial/Industrial applications

Facebook, Amazon, Microsoft, Google and Apple — are already offering this feature on various devices through services like Google Home, Amazon Echo and Siri.

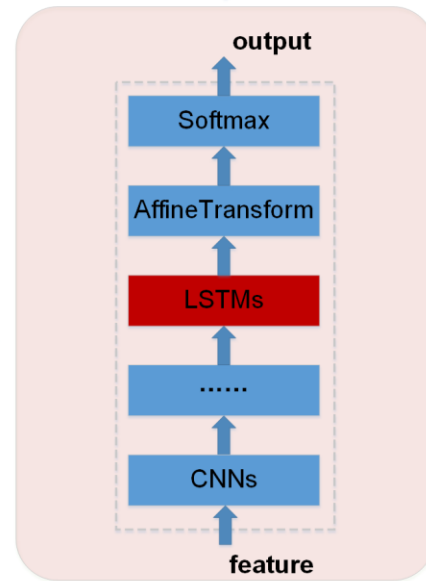
The try to **make speech recognition a standard for most product.**

# Speech to Text Application

- AI-based Solutions (LSTM)



$$\begin{aligned}
 i_t &= g(W_{ix}x_t + W_{ih}h_{t-1} + b_i) \\
 f_t &= g(W_{fx}x_t + W_{fh}h_{t-1} + b_f) \\
 o_t &= g(W_{ox}x_t + W_{oh}h_{t-1} + b_o) \\
 c\_in_t &= \tanh(W_{cx}x_t + W_{ch}h_{t-1} + b_c) \\
 c_t &= f_t \cdot c_{t-1} + i_t \cdot c\_in_t \\
 h_t &= o_t \cdot \tanh(c_t)
 \end{aligned}$$



**Abilities:**

(1) We have a demo base on Librispeech 1000h dataset under Deepspeech2 framework, which could show the entire flow of our algorithm, software and hardware co-design (containing pruning, quantization, compilation and FPGA inference);

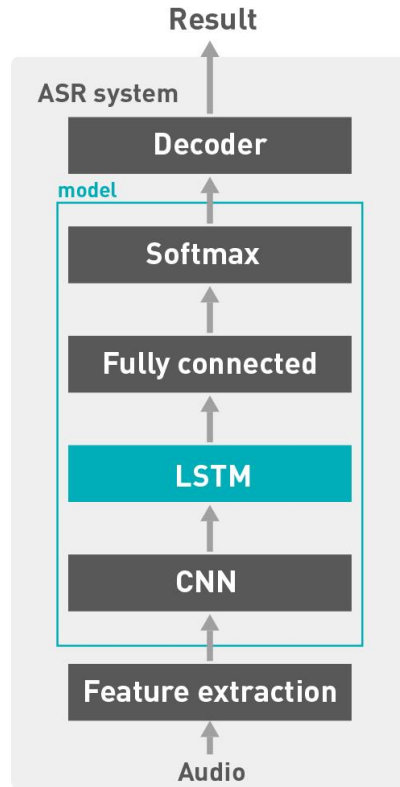
**Notes:**

(1) Compared to the traditional ASR system, a neural network with the same high accuracy is much more difficult to obtain. It means we haven't got a baseline with high accuracy.

## Applications in end-to-end ASR

# Speech to Text Application

- AI-based Solutions (LSTM)



1 layer only used in training

1 layer

5 layers Bi-directional

2 layers with BN and Hardtanh

Time (ms)	Only CPU on AWS	CPU + FPGA on AWS	CPU + GPU P4 locally (cudnn)
Decoding	0.09	0.08	0.40
Softmax	0.07	0.08	0.08
Fully connected	0.11	0.73	0.12
<b>LSTM</b>	<b>118.31</b>	<b>16.97</b>	<b>38.58</b>
CNN	14.16		1.21
Feature extraction	2.78	2.73	1.95
<i>E2E time</i>	<i>135.61</i>	<b>20.59</b> ★	<i>42.35</i>

**2.06X speedup!**

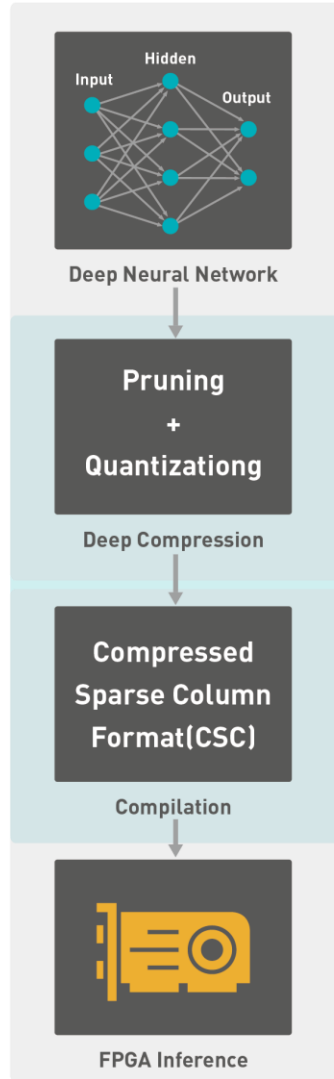
solutions	Devices and versions
CPU + FPGA on AWS	CPU: Intel(R) Xeon(R) CPU E5-2686 v4 @ 2.30GHz (8 processors) FPGA: VU9P
Only CPU on AWS	CPU: Intel(R) Xeon(R) CPU E5-2686 v4 @ 2.30GHz (8 processors)
CPU + GPU(P4) locally	CPU: Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz (56 processors) cuda: 8.0.44 cudnn: 6020

# LSTM Solution Overview



# Overview of LSTM solution

Flowchart

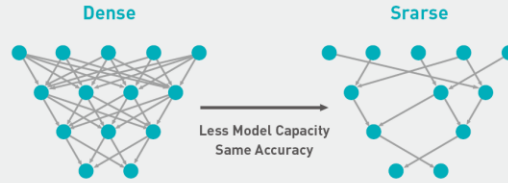


Deep Compression  
(Best paper of ICLR2016)

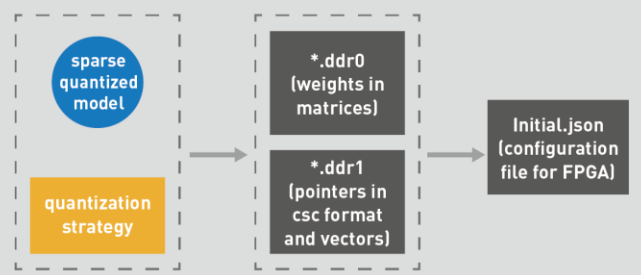
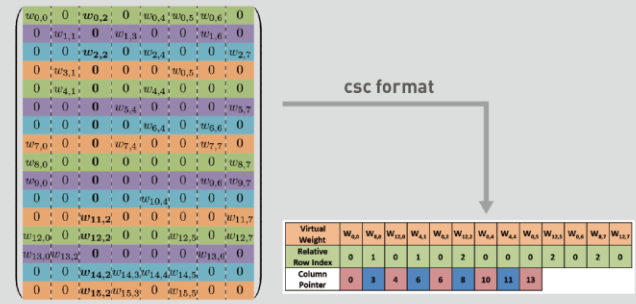
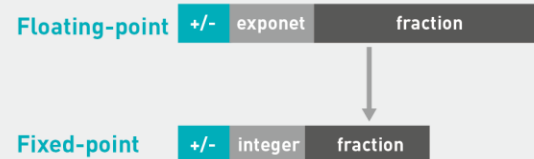
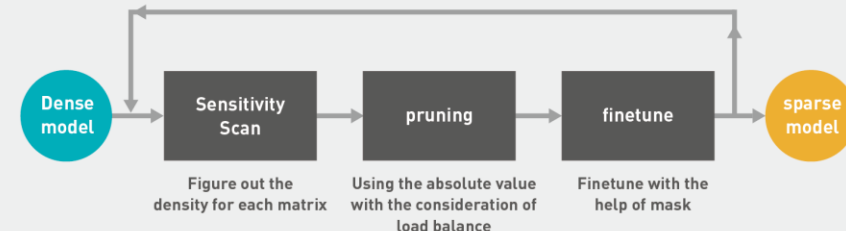
Software +

ESE  
(Best paper of FPGA2017)

Example



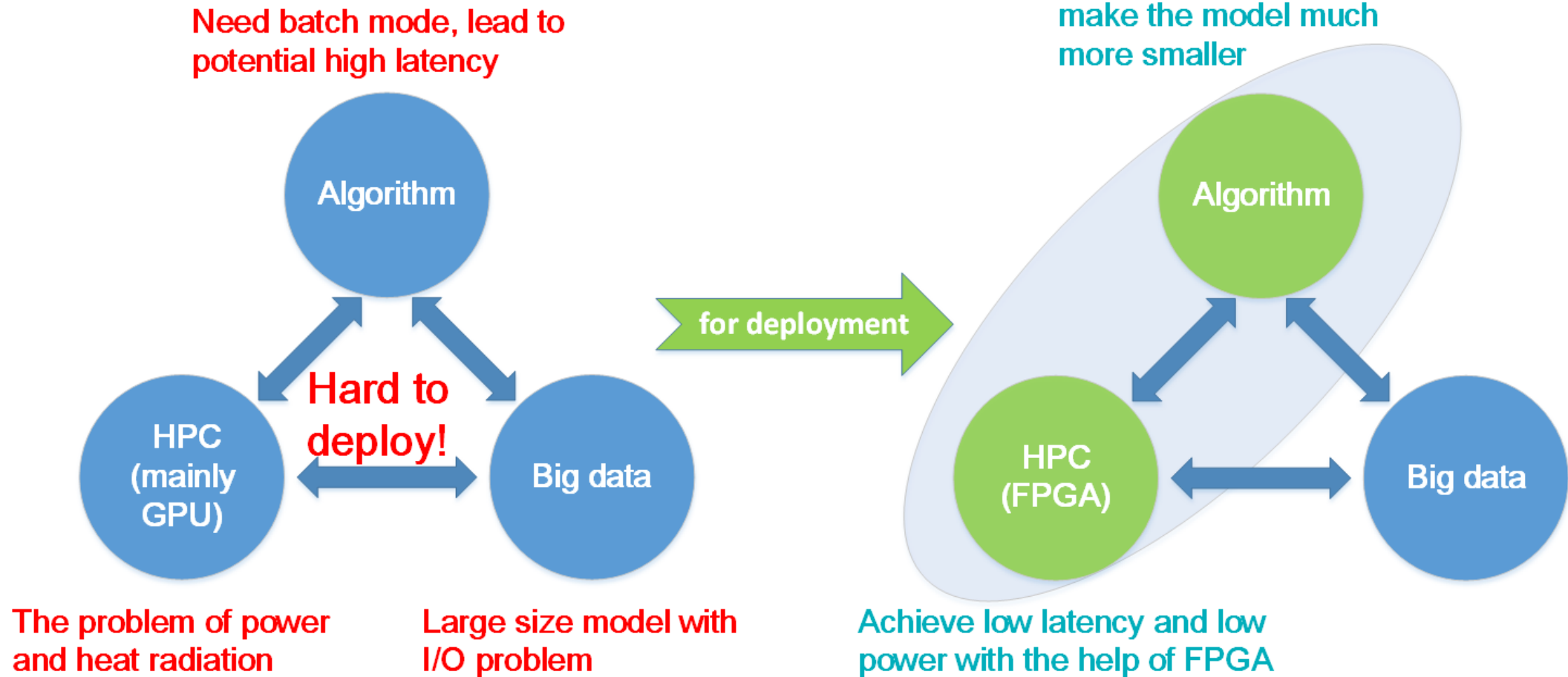
Algorithm Flow



Our algorithm, software and hardware co-design flow for RNN(LSTM) acceleration

# Overview of LSTM solution

What is the role of DeePhi in this game for LSTMs?



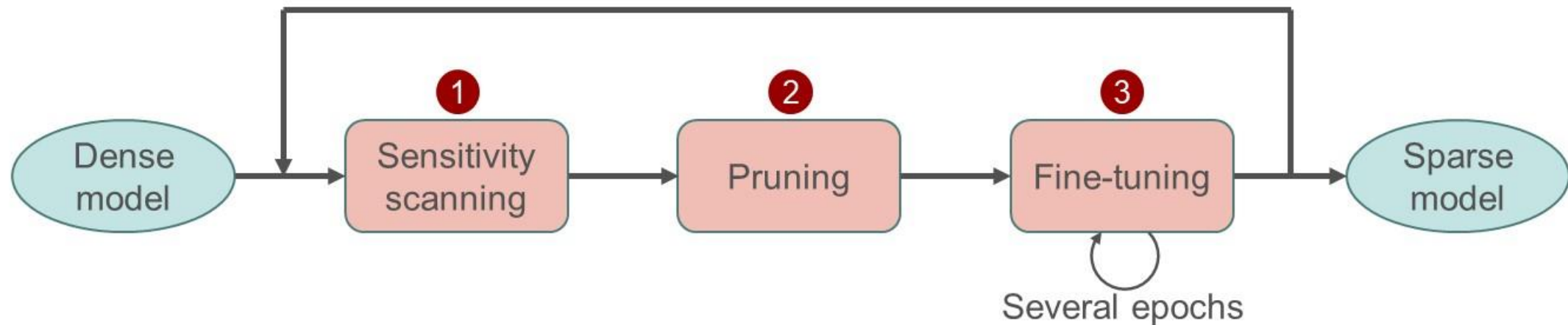


# Deep Compression for LSTM



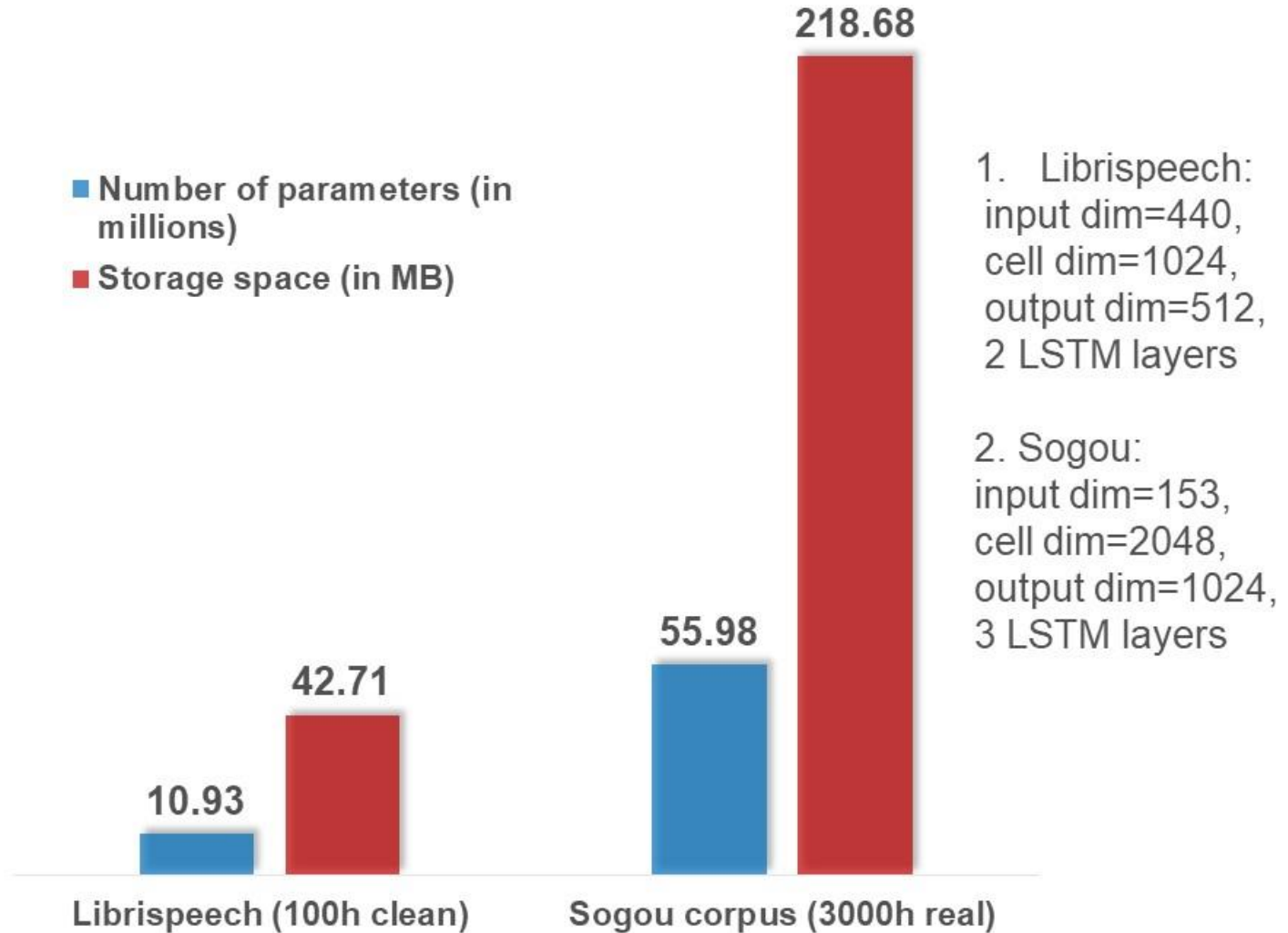
# Deep Compression for LSTM

- Step 1: Perform sensitivity scanning to determine how much to be pruned
- Step 2: Prune parameters based on absolute values
- Step 3: Fine-tune to enhance rest parameters
- Repeat step 1-3 to achieve higher sparsity



Sensitivity-based pruning flow for LSTM

# Deep Compression for LSTM



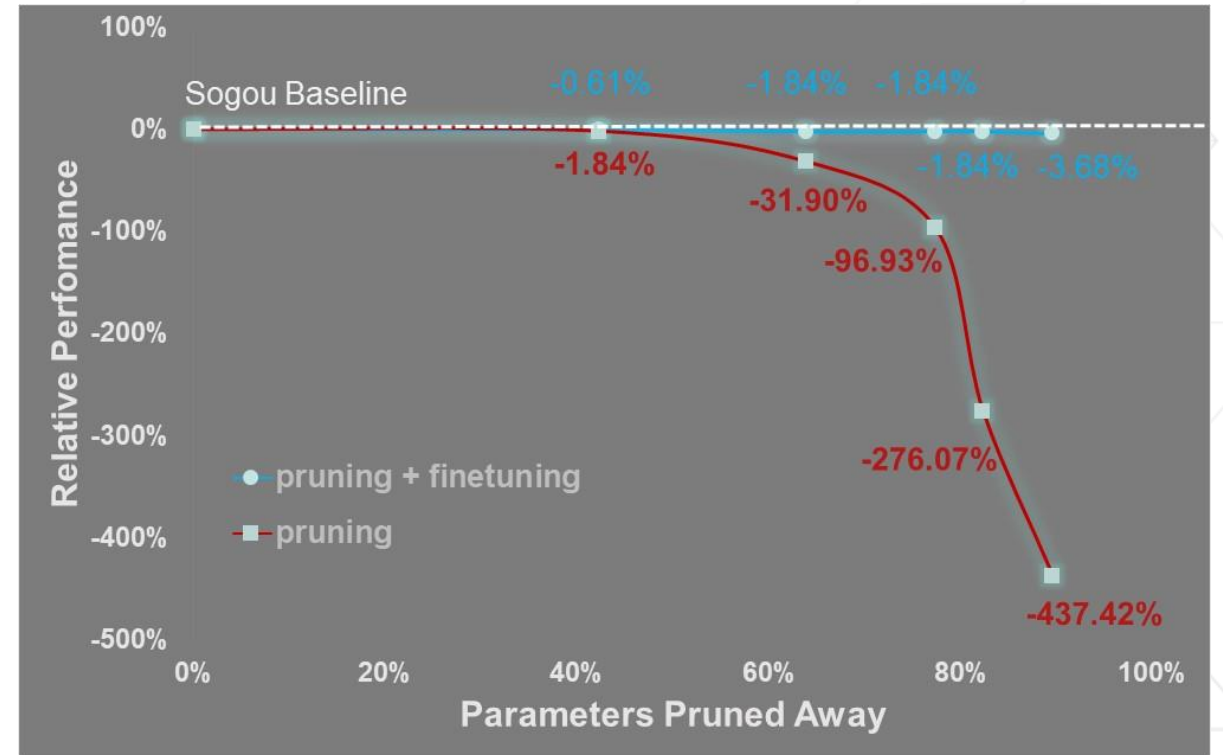
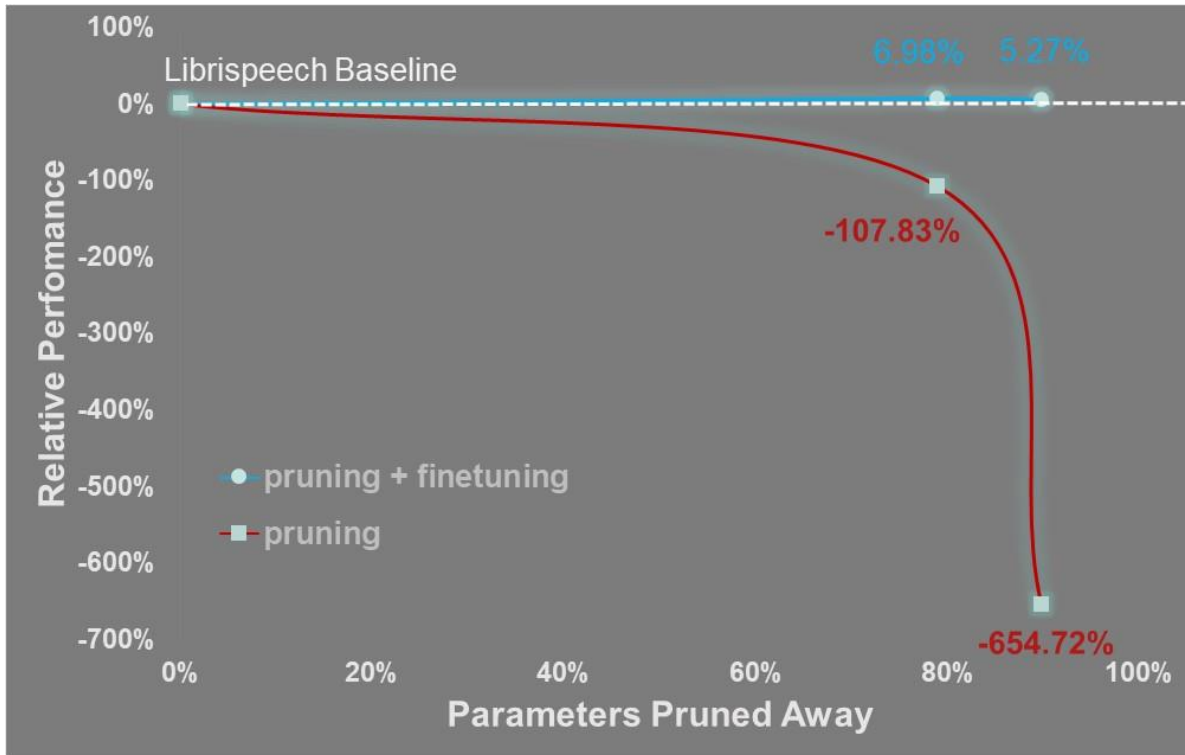
Examples

# Deep Compression for LSTM

Dataset	Model	density	WER	Relative Performance
Small (100h)	Librispeech dense	100%	12.90	0.0%
	Librispeech sparse	21.1%	12.00	+6.98%
	Librispeech sparse	10.30%	12.22	+5.27%
Big (3000h)	Sogou dense	100%	16.3	0.0%
	Sogou sparse	17.76%	16.6	-1.84%
	Sogou sparse	10.47%	16.9	-3.68%

Density vs accuracy

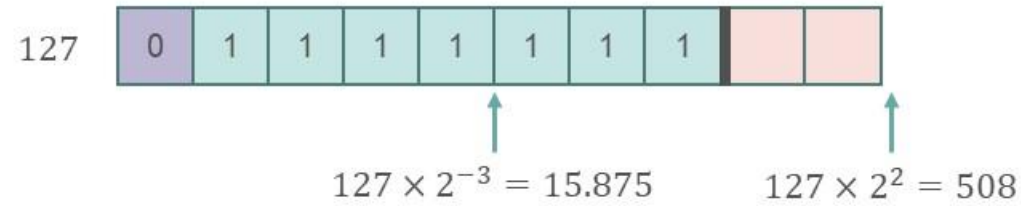
# Deep Compression for LSTM



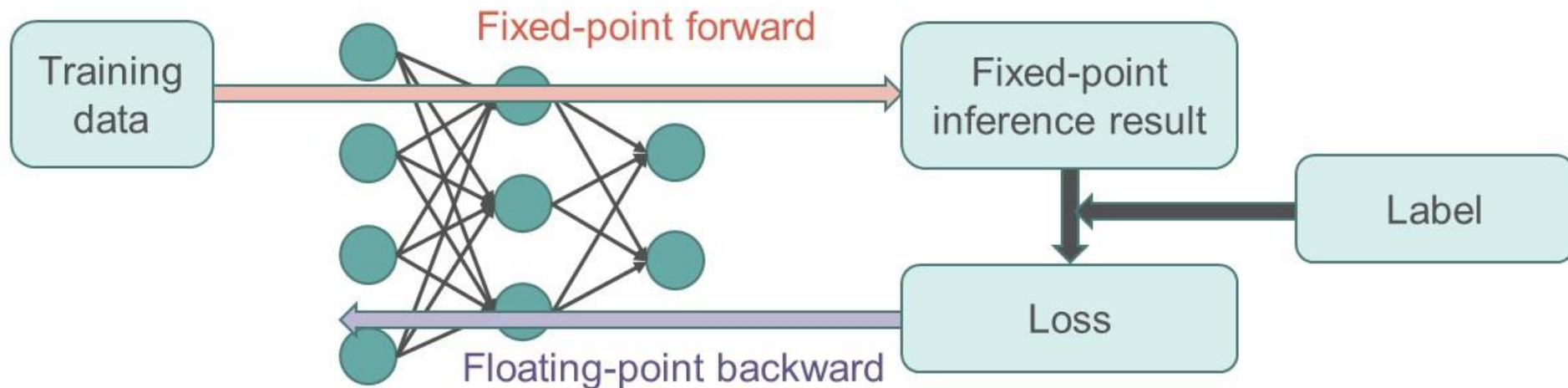
Sparsity vs accuracy

# Deep Compression for LSTM

## Quantization strategy (e.g. 8bit)



## Fixed-point re-training



## Quantization and fixed-training for LSTM

# Deep Compression for LSTM

Dataset	Model	density	WER	Relative Performance
Small (100h) English	Librispeech dense	100%	12.90	0.0%
	Librispeech sparse	10.30%	12.22	+6.98%
	Librispeech sparse and quantized	10.30%	11.96	+7.29%
Big (3000h) Chinese	Sogou dense	100%	16.3	0.0%
	Sogou sparse	17.76%	16.6	-1.84%
	Sogou sparse and quantized	17.76%	16.6	-1.84%

1. Quantization strategy:  
for Librispeech, weights and x/y/m vectors are both directly quantized to 8bits.  
for Sogou network, weights are quantized to 12bits and all vectors are quantized to 16bits;
2. WER: word error rate

Deep compression vs accuracy

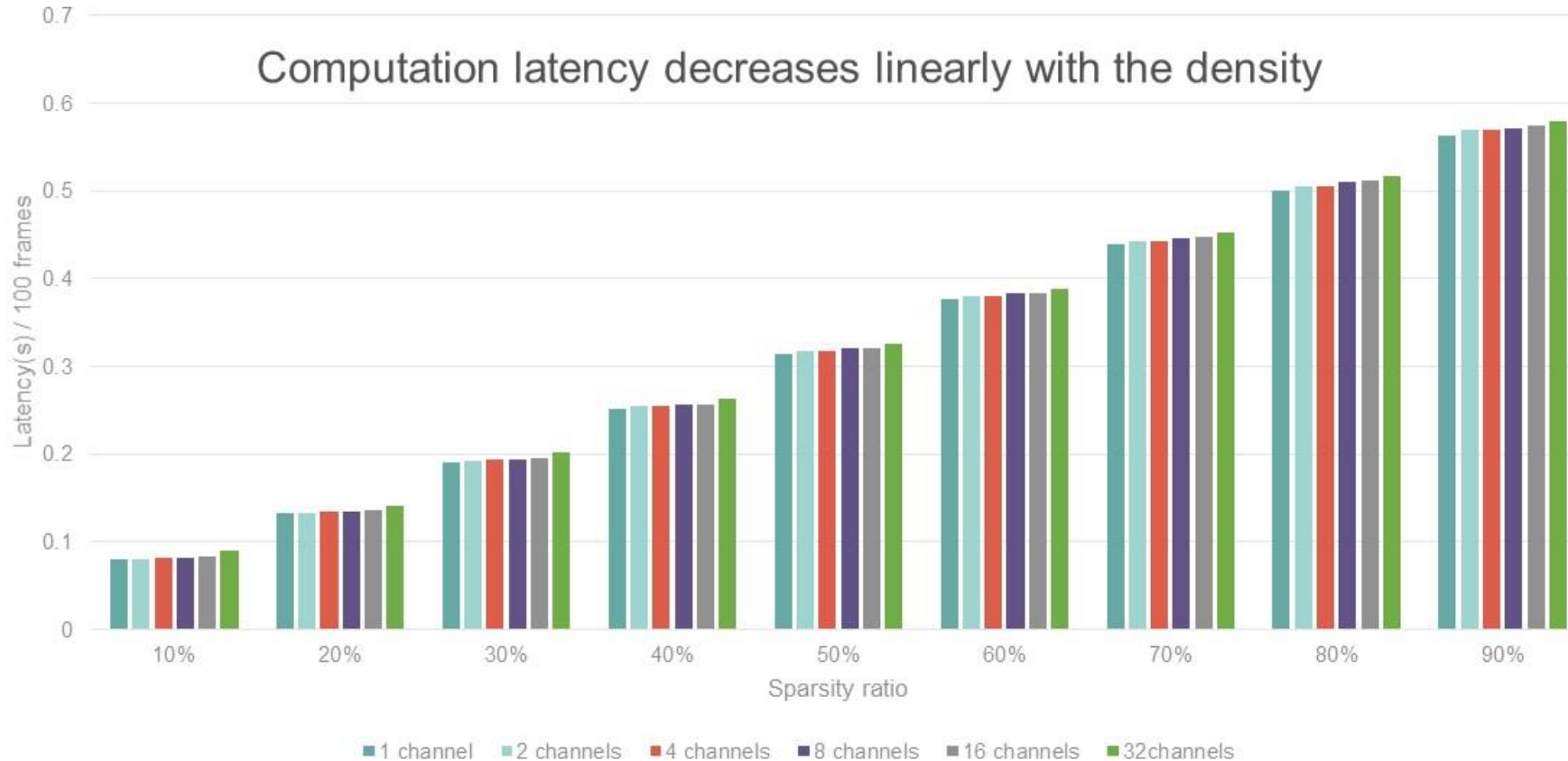
# Compressed LSTM on ESE





# Compressed LSTM on ESE

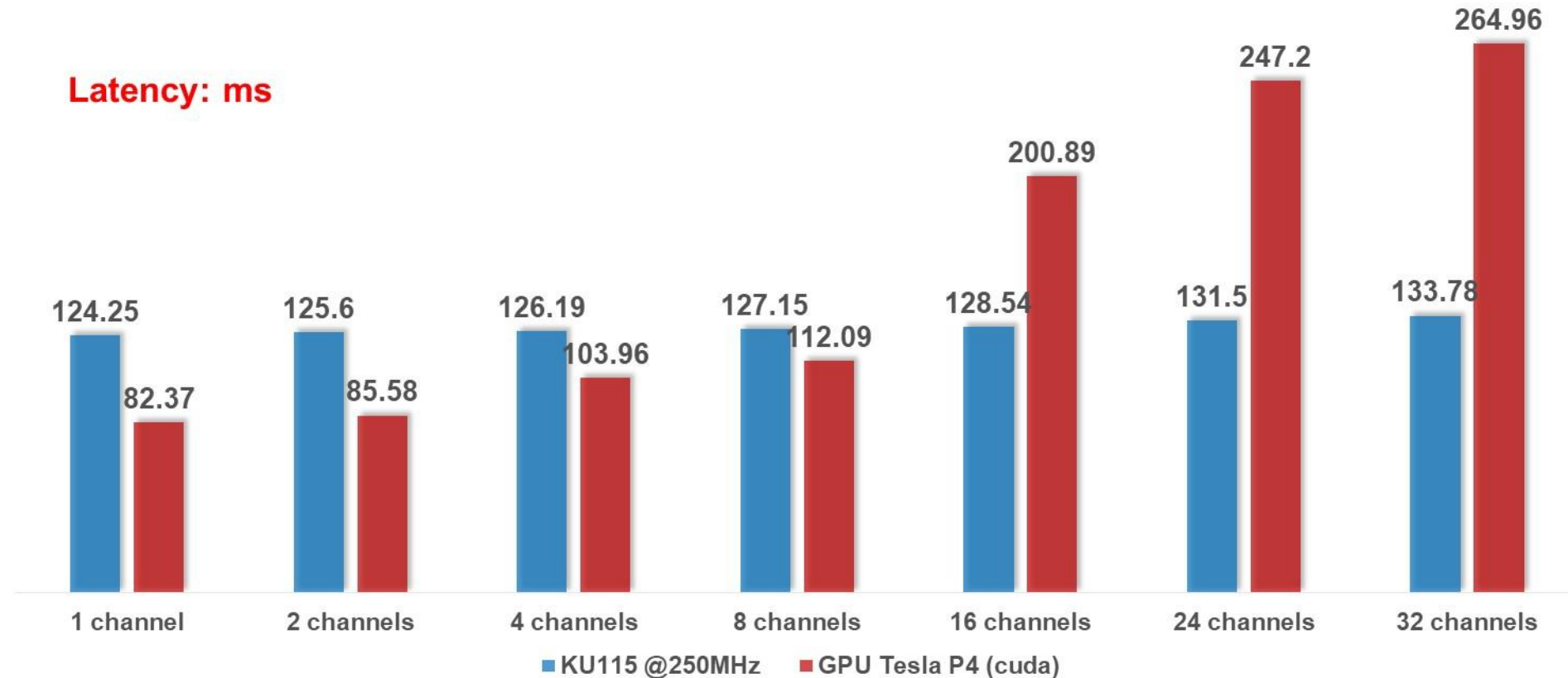
## Three-layers LSTM



1. Net parameters: input dim=153, cell dim=2048, output dim=1024, 3 LSTM layers;
2. KU115 is clocked at 250MHz;
3. CPU: Intel(R) Xeon(R) CPU E5-2690 v4 @2.60GHz 14 Cores 56 Threads

Latency under different density

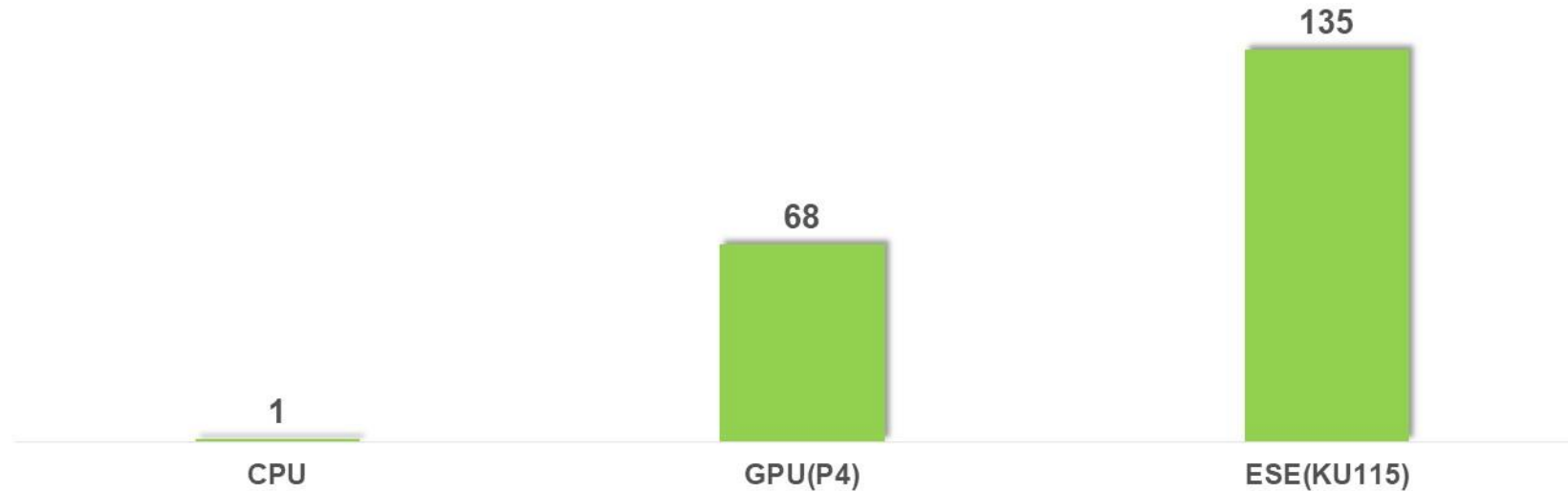
# Compressed LSTM on ESE



Parameters: input dim=153, cell dim=2048, output dim=1024, 3 LSTM layers, density=17.76%;  
Local CPU: Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz (56 processors)  
Non-merged version, 100 frames of input

Performance of LSTMP model for different channels based on Sogou 3000h dataset

# Compressed LSTM on ESE



Platforms	CPU(E5-2690 v4)	GPU(P4)	ESE(KU115)
Latency	18155.22ms	264.9576ms	133.778ms
Power	--	62 W	53 W
Speedup	1x	68x	135x

1. Nnet parameters: input dim=153, cell dim=2048, output dim=1024, 3 LSTM layers, density=17.76%;
2. ESE-KU115 is clocked at 250MHz;
3. 32 channel speech inputs in parallel.

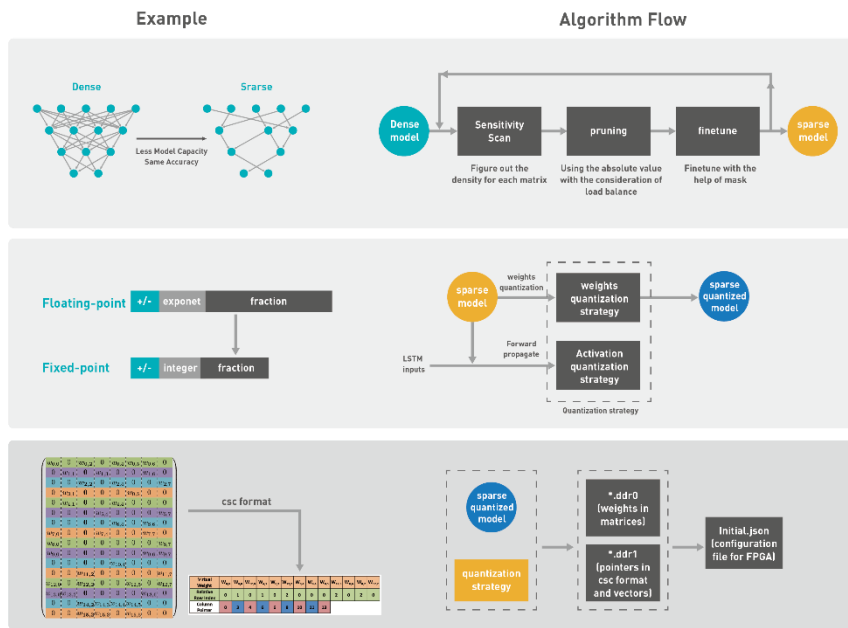
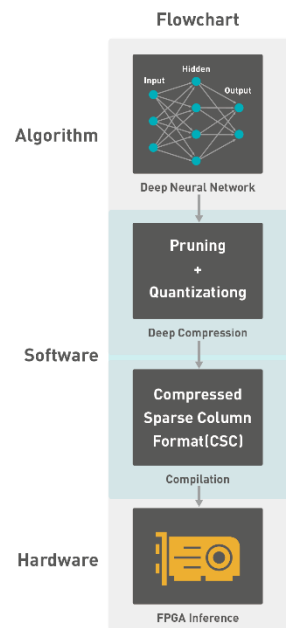
Performance of LSTMP model for 32 channels compared to CPU and GPU

# DDESE on Cloud

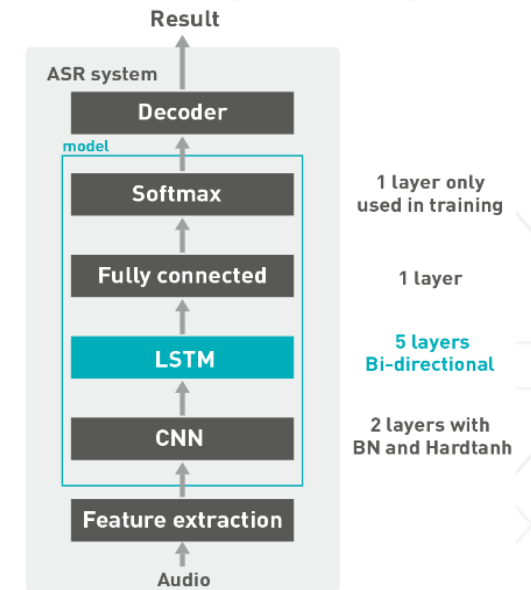


# DDESE on Clouds

## Solution



for end-to-end  
speech recognition



## Partners



✓ On clouds, aiming at customers all over the world



✓ Already officially launched in AWS Marketplace and HUAWEI cloud

(<http://www.deephi.com/ddese.html>)



✓ Now transplanting to Alibaba cloud

## Features

Low storage	Model compressed more than 10X with negligible loss of accuracy
Low latency	More than 2X speedup compared to GPU (P4)
Programmable	Reconfigurable for different requirements

# DDESE on Clouds

## The same model

### Algorithm

Accelerating CNN/BLSTM  
Pruning BLSTM to 15%  
16bit weights/activations

Accelerating CNN/BLSTM  
Pruning BLSTM to 15%  
16bit weights/activations

Accelerating LSTMP  
Pruning LSTMP to 26.9%  
16bit weights/activations

### Software

DeepSpeech2 + PyTorch  
Tools for compression and  
compilation

DeepSpeech2 + PyTorch  
Tools for compression and  
compilation

DeepSpeech2 + PyTorch  
Tools for compression and  
compilation

### Hardware

Based on VU9P  
220MHz  
1 channel

Based on VU9P  
200MHz  
1 channel

Based on KU115  
300MHz  
1 channel



# DDESE on Clouds

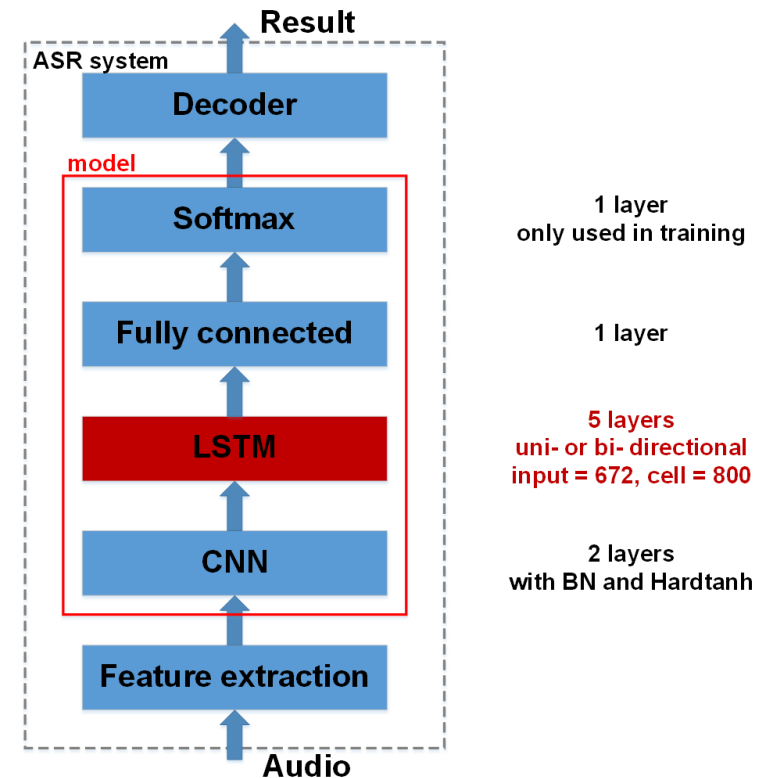
## Features

Innovative full-stack acceleration solution for deep learning in acoustic speech recognition (ESE: best paper of FPGA2017)

- Support both unidirectional and bi-directional LSTM acceleration on FPGA for model inference
- Support CNN layers, Fully-Connected (FC) layers, Batch Normalization layers and varieties of activation functions such as Sigmoid, Tanh and HardTanh
- Support testing for both performance comparison of CPU/FPGA and single sentence recognition
- Supporting user's own test audio recognition (English, 16kHz sample rate, no longer than 3 seconds)

Usage	Hardware PCIE interface, software API
Supported layer	CNN, uni/bi-directional LSTM(P), FC, BN
LSTM layer number	According to the requirements and source
Channel number	According to the requirements and source
Quantization	16bit
Maximum input of LSTM	1024
Maximum size of LSTM	2048
Density of LSTM	Any, typically 10%~20%
Peephole in LSTM	Selectable
Projection in LSTM	Selectable
Activation function	Sigmoid, Tanh, HardTanh

Note: It is available to change the hardware configuration for different requirements such as layer number and model size



# DDESE on AWS



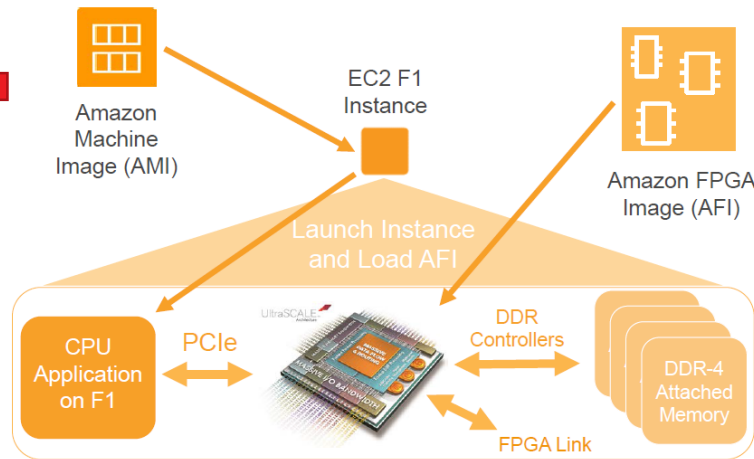
✓ Already officially launched in AWS Marketplace

[https://aws.amazon.com/marketplace/pp/B079N2J42R?qid=1523443241195&sr=0-1&ref\\_=srh\\_res\\_product\\_title](https://aws.amazon.com/marketplace/pp/B079N2J42R?qid=1523443241195&sr=0-1&ref_=srh_res_product_title)

	Version	Launch	Description
DDESE	V1.0	2017.12	acceleration for unidirectional and bi-directional LSTM model
	V2.0	2018.02	acceleration for CNN + bi-directional LSTM model

## Add DeePhi solution

### FPGA Acceleration Using F1



An F1 instance can have any number of AFIs  
An AFI can be loaded into the FPGA in less than 1 second



## Model Description

	Released model from GitHub	Our model (low accuracy)	Our model (high accuracy)
RNN type	GRU	LSTM(without projection)	LSTM(without projection)
RNN layers	5	5	5
RNN input	672	672	672
RNN size	800	800	800
bi-directional	yes	no	yes
batch normalization	yes	no	yes
WER on libri-test-clean	11.20	15.608	10.764
Best model after compressing LSTM	---	20% density, LSTM 16bit WER=17.033	15% density, LSTM 16bit WER=11.51
Best model after further quantizing CNN	---	---	CNN 16bit weight/ACT WER=11.52

David Pellerin from AWS in HotChips2017

>> 30



# Performance of DDESE on AWS

## ➤ For LSTM layers only:

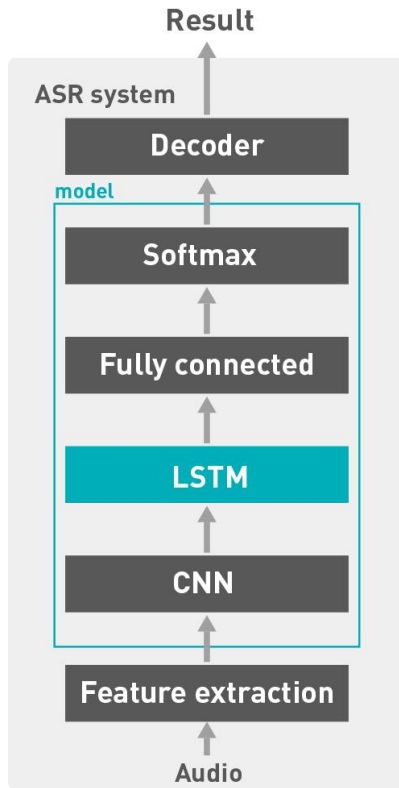
	Baseline (cudnn)	Compression	Hardware	Latency	Speedup
<i>Unidirectional LSTM</i>	WER = 15.608 E2E = 26.42ms <b>LSTM = 22.60ms</b>	WER = 17.033 Density = 20% 16bit weight/ACT	290MHz	E2E = 23.76ms <b>LSTM = 7.88ms</b>	1.11X <b>2.87X</b>
<i>Bi-directional LSTM</i>	WER = 10.764 E2E = 42.35ms <b>LSTM = 38.58ms</b>	WER = 11.51 Density = 15% 16bit weight/ACT	200MHz	E2E = 32.61ms <b>LSTM = 15.07ms</b>	1.30X <b>2.56X</b>

## ➤ For CNN + bi-directional LSTM layers:

	Baseline (cudnn)	Compression	Hardware	Latency	Speedup
<i>CNN+ BLSTM</i>	WER = 10.764 E2E = 42.35ms <b>CNN = 1.21ms</b> <b>LSTM = 38.58ms</b>	WER = 11.52 LSTM density = 15% All 16bit weight/ACT	<b>220MHz</b> (on going)	E2E 20.59ms <b>CNN+STM</b> <b>16.97ms</b>	2.06X

Note: E2E is short for end-to-end, ACT is short for activation, WER is short for word error rate, input: 1 second.

# Details of CNN + bi-directional model on AWS



1 layer only used in training

1 layer

5 layers Bi-directional

2 layers with BN and Hardtanh

	Time (ms)	Only CPU on AWS	CPU + FPGA on AWS	CPU + GPU P4 locally (cudnn)
Decoding		0.09	0.08	0.40
Softmax		0.07	0.08	0.08
Fully connected		0.11	0.73	0.12
LSTM		118.31	16.97	38.58
CNN		14.16		1.21
Feature extraction		2.78	2.73	1.95
<i>E2E time</i>		<i>135.61</i>	<b>20.59</b> ★	<i>42.35</i>

**2.06X speedup!**

solutions	Devices and versions
CPU + FPGA on AWS	CPU: Intel(R) Xeon(R) CPU E5-2686 v4 @ 2.30GHz (8 processors) FPGA: VU9P
Only CPU on AWS	CPU: Intel(R) Xeon(R) CPU E5-2686 v4 @ 2.30GHz (8 processors)
CPU + GPU(P4) locally	CPU: Intel(R) Xeon(R) CPU E5-2690 v4 @ 2.60GHz (56 processors) cuda: 8.0.44 cudnn: 6020

# DDESE on HUAWEI Cloud



- ✓ **Already officially launched on HUAWEI Cloud**
  - <https://app.huaweicloud.com/product/00301-110982-0--0>
- ✓ **Based on VU9P @200MHz, 1 channel**
- ✓ **Using CNN + bi-directional LSTM model (with high accuracy)**

## Model Description

### Using this model

	Released model from GitHub	Our model (low accuracy)	Our model (high accuracy)
RNN type	GRU	LSTM(without projection)	LSTM(without projection)
RNN layers	5	5	5
RNN input	672	672	672
RNN size	800	800	800
bi-directional	yes	no	yes
batch normalization	yes	no	yes
WER on libri-test-clean	11.20	15.608	10.764
Best model after compressing LSTM	---	20% density, LSTM 16bit WER=17.033	15% density, LSTM 16bit WER=11.51
Best model after further quantizing CNN	---	---	CNN 16bit weight/ACT WER=11.52

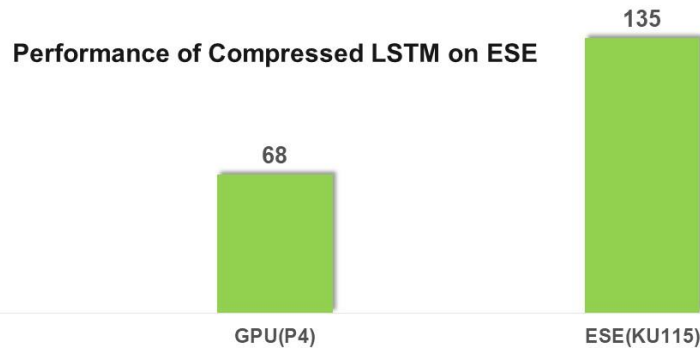
## Performance of CNN+BLSTM

Solution	Latency (ms)	Speedup
GPU (P4)	39.79	1
AWS @220MHz	16.97	2.34
HUAWEI @200MHz	18.60	2.14

# Review of Typical Products

Deep Compression vs Accuracy

Dataset	Model	density	WER	Relative Performance
Chinese (3000h real corpus)	Sogou dense	100%	16.3	0.0%
	Sogou sparse	17.76%	16.6	-1.84%
	Sogou sparse and quantized	17.76%	16.6	-1.84%



Platforms	CPU(E5-2690 v4)	GPU(P4)	ESE(KU115)
Latency	18155.22ms	264.9576ms	133.778ms
Power	--	62 W	53 W
Speedup	1x	68x	135x

1. Nnet parameters: input dim=153, cell dim=2048, output dim=1024, 3 LSTM layers, density=17.76%;
2. ESE-KU115 is clocked at 250MHz;
3. 32 channel speech inputs in parallel.

ESE for Sogou

## Partners



- ✓ On clouds, aiming at customers all over the world
- ✓ Already officially launched in AWS Marketplace and HUAWEI cloud (<http://www.deephi.com/ddese.html>)
- ✓ Soon available on Alibaba cloud

## Features

Low storage	Model compressed more than 10X with negligible loss of accuracy
Low latency	More than 2X speedup compared to GPU (P4)
Programmable	Reconfigurable for different requirements
High performance	Corresponding to 2.52 TOPS for dense model

DDESE on Clouds

# AI Boosting On-premise Solutions

DEEPhi



XILINX<sup>®</sup>  
ALVEO<sup>™</sup>

**Adaptable.**  
**Intelligent.**

