### Building Speech Recognition Systems with the Kaldi Toolkit

Sanjeev Khudanpur, Dan Povey and Jan Trmal Johns Hopkins University Center for Language and Speech Processing

June 13, 2016







#### In the beginning, there was nothing

• Then Kaldi was born in Baltimore, MD, in 2009.





#### Kaldi then grew up & became ...

← → C	he New York Times - Brei × Wp Washington Post: Breaking							
'KALDI'								
Main Page Related Pages Module	es Namespaces Classes Files							
▼ 'KALDI'	Kaldi							
<ul> <li>About the Kaldi project</li> <li>The build process (how Kaldi is compiled)</li> <li>Clustering mechanisms in Kaldi Contacting the Kaldi team</li> </ul>	Please see the instructions on upgrading your re Sourceforge. (see also Kaldi's project page on Sou • About the Kaldi project	epository to the new location, following o <b>60+ Contributors</b> urceforge ) Icon from http://thumbs.gograph.com						
<ul> <li>The CUDA Matrix library</li> <li>Data preparation</li> <li>Decoders used in the Kaldi toolkit</li> <li>Software required to install and run Kaldi</li> </ul>	<ul> <li>Other Kaldi-related resources</li> <li>Downloading and installing Kaldi</li> <li>Software required to install and run Kald</li> <li>Legal stuff</li> <li>Recent progress and current activity</li> </ul>	Postings to Discussion List						
<ul> <li>Deep Neural Networks in Kaldi</li> <li>Karel's DNN training implementation</li> <li>Deep Neural Networks in Kaldi (Dan's setup</li> <li>Kaldi logging and error-reporting</li> <li>Feature extraction</li> </ul>	<ul> <li>Kaldi tutorial</li> <li>Data preparation</li> <li>The build process (how Kaldi is compile</li> <li>The Kaldi coding style</li> <li>Contacting the Kaldi team</li> <li>History of the Kaldi project</li> </ul>	250 200 150						
<ul> <li>Finite State Transducer algorithms</li> <li>Decoding graph construction in Kaldi</li> <li>Decoding-graph creation recipe (test time)</li> <li>Decoding-graph creation recipe (training ti</li> <li>History of the Kaldi project</li> </ul>	<ul> <li>The Kaldi Matrix library</li> <li>External matrix libraries</li> <li>Kaldi I/O mechanisms</li> <li>Kaldi I/O from a command-line perspecti</li> <li>Kaldi logging and error-reporting</li> </ul>	100 50						
<ul> <li>HMM topology and transition modeling</li> <li>Downloading and installing Kaldi</li> <li>Kaldi I/O mechanisms</li> <li>Kaldi I/O from a command-line perspective</li> <li>Keyword Search in Kaldi</li> </ul>	<ul> <li>Parsing command-line options</li> <li>Other Kaldi utilities</li> <li>Clustering mechanisms in Kaldi</li> <li>HMM topology and transition modeling</li> <li>Decision tree internals</li> <li>How decision trees are used in Kaldi</li> </ul>	0 Mar-12 May-12 Jul-12 Sep-12 Jul-13 Jul-13 Jul-13 Sep-13 Jul-13 Jul-13 Jul-14 Jan-14 May-14 May-14 May-14						
	Decoding graph construction in Kaldi	Generated on Thu Jun 5 2014 12:46:52 for WAI DI by G DYN/G Dia 1 8 1 2						

#### Meanwhile, Speech Search went from "Solved" to "Unsolved" ... Again

- NIST TREC SDR (1998)
  - Spoken "document" retrieval from STT output as good as retrieval from reference transcripts
  - Speech search was declared a solved problem!
- NIST STD Pilot (2006)
  - STT was found to be inadequate for spoken "term" detection in conversational telephone speech
- Limited language diversity in CTS corpora
  - English Switchboard, Call Home and Fisher
  - Arabic and Mandarin Chinese Call Home

### In 2012, IARPA launched BABEL

One month after Dan Povey returned to Kaldi's birthplace

- Automatic transcription of conversational telephone speech was still the core challenge.
- But with a few subtle, crucial changes
  - Focused attention on low-resource conditions
  - Required concurrent progress in multiple languages
    - **PY1**: Cantonese, Tagalog, Pashto, Turkish and Vietnamese
    - PY2: Assamese, Bengali, Haitian Creole, Lao, Zulu and Tamil
  - Reduced system development time from year to year
  - Used keyword search metrics to measure progress

#### Kaldi Today

A community of Researchers Cooperatively Advancing STT

- C++ library, command-line tools, STT "recipes"
   Freely available via GitHub (Apache 2.0 license)
- Top STT performance in open benchmark tests

   E.g. NIST OpenKWS (2014) and IARPA ASpIRE (2015)
- Widely adopted in academia and *industry* 
  - 300+ citations in 2014 (based on Google scholar data)
  - 400+ citations in 2015 (based on Google scholar data)
  - Used by several US and non-US companies
- Main "trunk" maintained by Johns Hopkins

   Forks contain specializations by JHU and others

### Co-Pl's, PhD Students and Sponsors

- Sanjeev Khudanpur
- Daniel Povey
- Jan Trmal
- Guoguo Chen
- Pegah Ghahremani
- Vimal Manohar
- Vijayaditya Peddinti
- Hainan Xu
- Xiaohui Zhang
- and several others



# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building
  - Acoustic model training
  - Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring
- Basic DNN system building
- Going beyond the basics



#### Setting up Paths, Queue Commands, ...





# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building
  - Acoustic model training
  - Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring
- Basic DNN system building
- Going beyond the basics



#### Preparing Acoustic Training Data

			Т	erminal			
• • •			Tei	rminal			
/Users/san	jeev/Desktop/0	Confere	ence Trav	vel/	SLΤι	J 2016/	/Tutorial/kaldi/egs/iban/s5/\
data/train:							
total used	in directory	1944 (	available	e 21	3882	2280	
drwxr-xr-x	11 sanjeev	staff	374	May	3	12:11	
drwxr-xr-x	16 sanjeev	staff	544	May	3	12:25	
drwxr-xr-x	8 sanjeev	staff	272	May	3	10:21	.backup
-rw-rr	1 sanjeev	staff	1081	May	3	10:23	cmvn.scp
-rw-rr	1 sanjeev	staff	204475	May	3	10:23	feats.scp
-rw-rr	1 sanjeev	staff	32044	May	3	10:14	spk2utt
drwxr-xr-x	18 sanjeev	staff	612	May	3	10:29	split16
-rw-rr	1 sanjeev	staff	418273	May	3	10:14	text
-rw-rr	1 sanjeev	staff	54436	May	3	12:11	utt2dur
-rw-rr	1 sanjeev	staff	53180	May	3	10:14	utt2spk
-rw-rr	1 sanjeev	staff	220697	May	3	10:14	wav.scp

#### data/train/text

#### Terminal

ibf_002_165	bagiiban miri fm bisi nyediaka tempat ke empat bengkah program gaw\
ai rambau mu	sim pengerami gawai ke taun tu
ibf_002_166	program gawai tu nyengkaum jaku pesan ari penulung menteri perengk\
a ke diguna	mensia mayuh datuk sylvester entri muran enggau
ibf_002_167	program anak mit begawai ba studio ke dikeluarka mary bajang
ibf_002_169	kepala bagi iban rtm miri encik simon suti madahka bala nembiak ke∖
masuk progr	am setengah jam nya datai ari sk lambir village sk beluru central e∖
nggau sk kel	apa sawit numor dua
ibf_002_171	program nya deka ditabur ba serata menua nengah waifm ba ari ti ke\
teru <mark>b</mark> ah gawa	i pukul sepuluh pagi
ibf_002_172	sebengkah agi program gawai nya berami gawai miri fm dua ribu dua ∖
belas pengel	ama empat puluh lima minit
ibf_002_174	program nya deka dikeluar ba ari ti kedua gawai pukul dua ngalih h\
ari	
ibf_002_175	nya naka kami naburka berita ari waifm berita nya tadi ditusun bev\
erly kaur se	reta disalin shirley limban salam satu malaysia
ibf_002_177	pukul tujuh pagi
ibf_002_179	diatu kami naburka berita ari waifm
-uu-:F1	text 4% L89 (Fundamental)
Undo!	

#### data/train/wav.scp

•	Terminal
ibf_002_001 wav	/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_001.\
ibf_002_002 wav	<pre>/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_002.\</pre>
ibf_002_003 wav	<pre>/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_003.\</pre>
ibf_002_004 wav	<pre>/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_004.\</pre>
ibf_002_005 wav	<pre>/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_005.\</pre>
ibf_002_006 wav	<pre>/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_006.\</pre>
ibf_002_007 wav	<pre>/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_007.\</pre>
ibf_002_008 wav	<pre>/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_008.\</pre>
ibf_002_010 wav	/home/ubuntu/kaldi/egs/iban/s5/corpus/data/wav/ibf_002/ibf_002_010.∖
-uu-:F1	<pre>wav.scp Top L1 (Fundamental)</pre>

#### data/train/(utt2spk|spk2utt)

•			Terminal			
ibf_002_392	ibf_002		<mark> </mark> ibf_002	ibf_002_001	ibf_002_002	ibf_002\$
ibf_002_395	ibf_002		libf_003	ibf_003_001	ibf_003_003	ibf_003\$
ibf_002_399	ibf_002		libf_004	ibf_004_001	ibf_004_002	ibf_004\$
ibf_002_402	ibf_002		libf_005	ibf_005_001	ibf_005_002	ibf_005\$
ibf_002_403	ibf_002		libf_006	ibf_006_001	ibf_006_002	ibf_006\$
ibf_002_408	ibf_002		libf_007	ibf_007_001	ibf_007_002	ibf_007\$
ibf_002_409	ibf_002		libf_008	ibf_008_001	ibf_008_002	ibf_008\$
ibf_002_410	ibf_002		libf_010	ibf_010_001	ibf_010_003	ibf_010\$
ibf_002_411	ibf_002		libf_014	ibf_014_001	ibf_014_002	ibf_014\$
ibf_003_001	ibf_003		libf_015	ibf_015_002	ibf_015_003	ibf_015\$
ibf_003_003	ibf_003		libm_001	ibm_001_001	ibm_001_002	ibm_001\$
ibf_003_004	ibf_003		libm_002	ibm_002_001	ibm_002_002	ibm_002\$
ibf_003_005	ibf_003		libm_003	ibm_003_002	ibm_003_003	ibm_003\$
ibf_003_009	ibf_003		libm_004	ibm_004_001	ibm_004_002	ibm_004\$
ibf_003_010	ibf_003		libm_006	ibm_006_001	ibm_006_002	ibm_006\$
ibf_003_011	ibf_003		libm_007	ibm_007_001	ibm_007_002	ibm_007\$
ibf_003_012	ibf_003		libm_010	ibm_010_001	ibm_010_003	ibm_010\$
ibf_003_013	ibf_003					
-uu-:F1	utt2spk	9% L229	(l-uu-:	-F1 spk2utt	All	L1 (F

### data/train/(cmvn.scp|feats.scp)

Terminal	
<pre>ibf_002 /home<sup>#</sup> initialization PATH ./path.sh II die "path.sh expected"; ibf_003 /home<sup>#</sup> initialization commands ./cmd.sh</pre>	
LDT_004 / NOME ./utils/parse_options.sh	
10 <sup>†</sup> _00 <sup>5</sup> / Nome set -e -o pipefail	
ibf_006 /home # download iban to build ASR	
ibf_007 /home <sup>if[!-f"\$corpus/README"];</sup>	
ibf_008 /home mkdir -p ./\$corpus/	
ibf_010 /home tar xzf iban.tar.gz -C \$corpus	ýz
-uu-:F1 c <sup>fi</sup>	
ibf_002_001 / <sup>nj=16</sup>	k:12
ibf_002_002 /	k:13877
ibf_002_003 / echo "Preparing data and training language models"	k:32877
ibf_002_004 / local/prepare_data.sh \$corpus/	k:38292
ibf 002 005 / utils/prepare_lang.sh data/local/dict " <unk>" data/local/lang data/lang</unk>	k:56876
ibf 002 006 / <sup>fi</sup>	k:71261
ibf 002 007 /	k:84346
if [ \$stage -le 2 ]; then ibf 002 008 / # Feature extraction	k·104425
ihf $002 010 / $ for x in train dev; do	k·121202
steps/make_mfcc.shnj \$njcma "\$train_cma" aata/\$x exp/make_mfcc/\$x mfc steps/compute_cmvn_stats.sh data/\$x exp/make_mfcc/\$x mfcc	
fi done	
1 = 1 $1 = 1$ $1 =$	

# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building
  - Acoustic model training
  - Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring
- Basic DNN system building
- Going beyond the basics



#### **Preparing the Pronunciation Lexicon**

$\bullet $				Te	ermina	al		
/Users/sanj	eev	//Desktop	/Confe	rence Tro	avel,	/SL1	FU 2016	6/Tutorial/kaldi/egs/iban/s5∖
/data/local/d	lict	::						
total used	in	director	y 3408	availab	le 23	1387	71436	
drwxr-xr-x	9	sanjeev	staff	306	May	2	20:07	
drwxr-xr-x	5	sanjeev	staff	170	May	2	20:15	<u>.</u> .
-rw-rr	1	sanjeev	staff	0	May	3	10:14	extra_questions.txt
-rw-rr	1	sanjeev	staff	788889	May	3	10:14	lexicon.txt
-rw-rr	1	sanjeev	staff	934289	May	2	20:07	lexiconp.txt
-rw-rr	1	sanjeev	staff	78	May	3	10:14	nonsilence_phones.txt
-rw-rr	1	sanjeev	staff	6	May	3	10:14	oov.txt
-rw-rr	1	sanjeev	staff	4	May	3	10:14	optional_silence.txt
-rw-rr	1	sanjeev	staff	4	May	3	10:14	<pre>silence_phones.txt</pre>



#### data/local/dict/lexicon.txt

•		Terminal
<sil></sil>	SIL	
<unk></unk>	SIL	
ke	k @	
nya	NJ a KK	
iya	ija	
ba	b a KK	
dua	duwa	
sida	sida KK	
puluh	pulu@h	
raban	raban	
lalu	lalu	
agi	agi KK	
orang	ura NG	
dot	dot	
ribu	ribu	
perintal	n p@rintah	
tiga	tiga	
menteri	m@ntri	
-uu-:	-F1 lexicon.txt Top L1	(Text)

#### data/local/dict/\*silence\*.txt



#### data/local/lang

• • •				Ter	minal			
/Users/sanj	eev/D	esktop/	Confer	ence Tra	vel/S	SLTU	J 2016/	/Tutorial/kaldi/egs/iban/s5 $lacksquare$
/data/local/l	ang:							
total used	in di	rectory	8064	availabl	e 213	3810	0860	
drwxr-xr-x	7 sa	njeev	staff	238	May	7	07:17	
drwxr-xr-x	5 sa	njeev	staff	170	May	2	20:15	<u>.</u> .
-rw-rr	1 sa	njeev	staff	1271927	May	3	10:14	align_lexicon.txt
-rw-rr	1 sa	njeev	staff	2	May	3	10:14	lex_ndisambig
-rw-rr	1 sa	njeev	staff	1417317	May	3	10:14	lexiconp.txt
-rw-rr	1 sa	njeev	staff	1425210	May	3	10:14	lexiconp_disambig.txt
-rw-rr	1 sa	njeev	staff	729	Мау	3	10:14	phone_map.txt

(Dired by name)---

#### Word Boundary Tags

•				Terminal
<sil> <unk> ke nya iya ba dua sida</unk></sil>	1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0	SIL SIL k@ NJ a KK i j a b a KK d u w a s i d a KK		
-uu-:	F1 lex	iconp.txt<2	> Top L1	(Text)
<sil> 1. <unk> 1. ke 1.0 k nya 1.0 iya 1.0 ba 1.0 k dua 1.0 sida 1.0 puluh 1.</unk></sil>	0 SIL_S 0 SIL_S (_B @_E NJ_B a_3 i_B j_I 0_B a_I 1 d_B u_I 0 s_B i_3 0 p_B u	I KK_E a_E KK_E w_I a_E I d_I a_I K _I l_I u_I	K_E @_I h_E	
-uu-:	F1 lex	iconp.txt	Top L1	(Text)

#### **Disambiguation Symbols**

•			Тег	minal				
<sil></sil>	1.0	SIL_S #1		lbailey	1.0	b_B e_I	l_I i_E	
<unk></unk>	1.0	SIL_S #2		lbailout	1.0	b_B e_I	i_I l_I	aw_I t_\$
ke	1.0	k_B @_E		lbain	1.0	b_B aj_]	[ n_E #1	
nya	1.0	NJ_B a_I K	K_E	lbait	1.0	b_B aj_]	[ t_E #1	
iya	1.0	i_B j_I a_	E	lbaja	1.0	b_B a_I	dZ_I @_E	
ba	1.0	b_B a_I KK	_E #1	lbak	1.0	b_B a_I	KK_E #2	
dua	1.0	d_B u_I w_	I a_E	lbaka	1.0	b_B a_I	k_I a_E	
sida	1.0	s_B i_I d_	I a_I KK_E	lbaker	1.0	b_B e_I	k_I @_E	
puluh	1.0	p_B u_I l_	I u_I @_I h_\$	lbakery	1.0	b_B e_I	k_I @_I	r_I i_E\$
raban	1.0	r_B a_I b_	I a_I n_E	lbaking	1.0	b_B a_I	k_I i_I I	NG_E
lalu	1.0	l_B a_I l_	I u_E	lbaku	1.0	b_B a_I	k_I u_E	
agi	1.0	a_B g_I i_	I KK_E	lbal	1.0	b_B a_I	l_E #1	
orang	1.0	u_B r_I a_	I NG_E #1	lbala	1.0	b_B a_I	l_I a_E	#1
dot	1.0	d_B o_I t_	E	lbalan	1.0	b_B a_I	l_I a_I	n_E
ribu	1.0	r_B i_I b_	I u_E	lbalas	1.0	b_B a_I	l_I a_I	s_E
perintał	า	1.0 p_	3 @_I r_I i_\$	lbaldi	1.0	b_B a_I	l_I d_I	i_E
tiga	1.0	t_B i_I g_	I a_E	lbale	1.0	b_B a_I	l_E #2	
menteri	1.0	m_B @_I n_	It_Ir_Ii_\$	lbali	1.0	b_B a_I	l_I i_E	
-uu-:	F1 lex	iconp_disam	big.txt Top	-uu-:	-F1 <b>lex</b>	iconp_dis	sambig.tx	t 89%

#### data/lang

• • •			Term	inal			
/Users/sanj	eev/Desktop/	Confere	ence Trave	el/SL	TU	2016/7	Tutorial∕kaldi/egs/iban/s5∖
/data/lang:							
total used	in directory	28128	available	e 213	813	3416	
drwxr-xr-x	10 sanjeev	staff	340	May	2	20:07	
drwxr-xr-x	16 sanjeev	staff	544	May	3	12:25	
-rw-rr	1 sanjeev	staff	6908222	May	3	10:14	L.fst
-rw-rr	1 sanjeev	staff	6981934	May	3	10:14	L_disambig.fst
-rw-rr	1 sanjeev	staff	2	May	3	10:14	oov.int
-rw-rr	1 sanjeev	staff	6	May	3	10:14	oov.txt
drwxr-xr-x	30 sanjeev	staff	1020	May	2	20:07	phones
-rw-rr	1 sanjeev	staff	1146	May	3	10:14	phones.txt
-rw-rr	1 sanjeev	staff	1369	May	3	10:14	topo
-rw-rr	1 sanjeev	staff	490107	May	3	10:14	words.txt

-uuu:%%-F1 lang<2>

#### data/lang/(phones|words).txt

			Terminal			
<eps> 0</eps>			I <mark>&lt;</mark> eps≻ 0			
SIL 1			I <unk> 1</unk>			
SIL_B 2			l <sil> 2</sil>			
SIL_E 3			la 3			
SIL_I 4			l <mark>a-lelaki 4</mark>			
SIL_S 5			l <mark>a-one 5</mark>			
@_B 6			l <mark>a-satu 6</mark>			
@_E 7			l <mark>aa</mark> 7			
@_I 8			l <mark>aabar 8</mark>			
@_S 9			l <mark>aad 9</mark>			
GG_B 10			l <mark>aadk 10</mark>			
GG_E 11			l <mark>aaj 11</mark>			
GG_I 12			l <mark>aaja 12</mark>			
GG_S 13			l <mark>aam 13</mark>			
KK_B 14			l <mark>aamir 14</mark>			
KK_E 15			l <mark>aank 15</mark>			
KK_I 16			laao 16			
KK_S 17			laari 17			
-uu-:F1	<pre>phones.txt</pre>	Top L1	(l-uu-:F1	words.txt	Top L1	(T

#### data/lang/topo

	Termina					
<forphones> 1 2 3 4 5 </forphones>						
<pre><state> 0 <pdfclass> 0 <transition> 0 .25 <transition> 3 0.25 </transition></transition></pdfclass></state></pre>	0 0.25	<transition></transition>	1 0.25	<transition></transition>	2	0\
<pre><state> 1 <pdfclass> 1 <transition> 1 .25 <transition> 4 0.25 </transition></transition></pdfclass></state></pre>	1 0.25	<transition></transition>	2 0.25	<transition></transition>	3	0\
<pre><state> 2 <pdfclass> 2 <transition> 3 .25 <transition> 4 0.25 </transition></transition></pdfclass></state></pre>	1 0.25	<transition></transition>	2 0.25	<transition></transition>	3	0\
<pre><state> 3 <pdfclass> 3 <transition> 1 .25 <transition> 4 0.25 </transition></transition></pdfclass></state></pre>	1 0.25	<transition></transition>	2 0.25	<transition></transition>	3	0\
<state> 4 <pdfclass> 4 <transition> 4 <state> 5 </state>  </transition></pdfclass></state>	4 0.75	<transition></transition>	5 0.25			

#### data/lang/phones/roots.txt

$\bullet$ $\bullet$		Terminal
shared	<pre>split SIL SIL_B SIL_E SIL_I</pre>	SIL_S
shared	split @_B @_E @_I @_S	
shared	split GG_B GG_E GG_I GG_S	
shared	split KK_B KK_E KK_I KK_S	
shared	split NG_B NG_E NG_I NG_S	
shared	split NJ_B NJ_E NJ_I NJ_S	
shared	<pre>split SS_B SS_E SS_I SS_S</pre>	
shared	split a_B a_E a_I a_S	
shared	split aj_B aj_E aj_I aj_S	
shared	split aw_B aw_E aw_I aw_S	
shared	split b_B b_E b_I b_S	
shared	split d_B d_E d_I d_S	
shared	<pre>split dZ_B dZ_E dZ_I dZ_S</pre>	
shared	split e_B e_E e_I e_S	
shared	split f_B f_E f_I f_S	
shared	split g_B g_E g_I g_S	
shared	split h_B h_E h_I h_S	
shared	split i_B i_E i_I i_S	
-uu-:	F1 roots.txt Top L1	(Text)

#### data/lang/phones/extra\_questions.txt

Terminal @\_B GG\_B KK\_B NG\_B NJ\_B SS\_B a\_B aj\_B aw\_B b\_B d\_B dZ\_B e\_B f\_B g\_B h\_B i\_B j\_B∖ k\_B l\_B m\_B n\_B o\_B oj\_B p\_B r\_B s\_B t\_B tS\_B u\_B v\_B w\_B x\_B z\_B @\_E GG\_E KK\_E NG\_E NJ\_E SS\_E a\_E aj\_E aw\_E b\_E d\_E dZ\_E e\_E f\_E g\_E h\_E i\_E j\_E k\_E l\_E m\_E n\_E o\_E oj\_E p\_E r\_E s\_E t\_E tS\_E u\_E v\_E w\_E x\_E z\_E <u>@\_I GG\_I KK\_I NG\_I NJ\_I SS\_I a\_I aj\_I aw\_I b\_I d\_I dZ\_I e\_I f\_I g\_I h\_I i\_I j\_I\</u> <u>k\_I l\_I m\_I n\_I o\_I oj\_I p\_I r\_I s\_I t\_I tS\_I u\_I v\_I w\_I x\_I z\_I</u> @\_S GG\_S KK\_S NG\_S NJ\_S SS\_S a\_S aj\_S aw\_S b\_S d\_S dZ\_S e\_S f\_S g\_S h\_S i\_S j\_S  $\langle$ k\_S l\_S m\_S n\_S o\_S oj\_S p\_S r\_S s\_S t\_S tS\_S u\_S v\_S w\_S x\_S z\_S SIL SIL\_B SIL\_E SIL\_I SIL S

# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building
  - Acoustic model training
  - Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring
- Basic DNN system building
- Going beyond the basics



#### Preparing the Language Model

 $\bullet$   $\bullet$ 

Terminal

local/train\_lms\_srilm.sh --train-text data/train/text data/ data/srilm

nl -nrz -w10 corpus/LM/iban-bp-2012.txt | sort -R > data/local/external\_text
local/train\_lms\_srilm.sh --train-text data/local/external\_text data/ data/srilm\
\_external

# let's do ngram interpolation of the previous two LMs # the lm.gz is always symlink to the model with the best perplexity, so we use  $\setminus$  that

#### local/train\_lms\_srilm.sh

•	Terminal		
echo	, ""		
echo echo	"Kneser-Ney 3grams" ""		
ngra nin cab	m-count -lm <mark>\$tgtdir</mark> /3gram.kn011.gz -kndisc 1 -kndiscount3 -gt3min 1 -order 3 -text <mark>\$t</mark> -unk -sort -map-unk "\$oov_symbol"	ount1 -gt1min 0 <mark>gtdir</mark> /train.txt	-kndiscount2 -gt2\ -vocab \$ <mark>tgtdir</mark> /vo\
ngra nin cab	m-count -lm \$ <mark>tgtdir</mark> /3gram.kn012.gz -kndisc 1 -kndiscount3 -gt3min 2 -order 3 -text \$ <mark>t</mark> -unk -sort -map-unk "\$oov_symbol"	ount1 -gt1min 0 <mark>gtdir</mark> /train.txt	-kndiscount2 -gt2\ -vocab \$ <mark>tgtdir</mark> /vo\
ngra nin cab	am-count -lm \$ <mark>tgtdir</mark> /3gram.kn022.gz -kndisc 2 -kndiscount3 -gt3min 2 -order 3 -text \$ <mark>t</mark> -unk -sort -map-unk "\$oov_symbol"	ount1 -gt1min 0 <mark>gtdir</mark> /train.txt	-kndiscount2 -gt2\ -vocab \$ <mark>tgtdir</mark> /vo\
ngra nin cab	am-count -lm \$ <mark>tgtdir</mark> /3gram.kn023.gz -kndisc 2 -kndiscount3 -gt3min 3 -order 3 -text \$ <mark>t</mark> -unk -sort -map-unk "\$ooy symbol"	ount1 -gt1min 0 gtdir/train.txt	-kndiscount2 -gt2\ -vocab \$tgtdir/vo\

#### local/train\_lms\_srilm.sh (cont'd)

• • •		Terminal	
-rw-rr-	- 1 sanjeev staf	F 33901333 May	3 10:19 4gram.me.gz
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 1240177 May	3 10:15 dev.txt
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 1376724 May	3 10:15 dev_text
lrwxr-xr-	x 1 sanjeev staf <sup>.</sup>	f 11 May	3 10:20 lm.gz -> 3gram.me.gz
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 5217 May	3 10:20 perplexities.txt
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 11103581 May	3 10:15 train.txt
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 12333505 May	3 10:15 train_text
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 283078 May	3 10:15 vocab
-uuu:%%-F1	srilm_external	Bot L37 (D	ired by name)
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 1332577 May	3 10:14 4gram.me.gz
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 38890 May	3 10:14 dev.txt
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 42070 May	3 10:14 dev_text
lrwxr-xr-	x 1 sanjeev staf <sup>.</sup>	f 11 May	3 10:15 lm.gz -> 4gram.me.gz
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 4639 May	3 10:15 perplexities.txt
-rw-rr-	- 1 sanjeev staf <sup>.</sup>	f 347475 May	3 10:14 train.txt
- rw-rr-	- 1 sanjeev staf	f 376203 May	3 10:14 train_text
-rw-rr-	- 1 sanjeev staf	F 283078 May	3 10:14 vocab

-uuu:%%-F1 srilm

Bot L41 (D

#### Interpolated Language Models

Terminal # let's do ngram interpolation of the previous two LMs # the lm.gz is always symlink to the model with the best perplexity, so we use  $\setminus$ that mkdir -p data/srilm\_interp for w in 0.9 0.8 0.7 0.6 0.5; do ngram -lm data/srilm/lm.gz -mix-lm data/srilm\_external/lm.gz \ -lambda \$w -write-lm data/srilm\_interp/lm.\${w}.gz echo -n "data/srilm\_interp/lm.\${w}.gz ' ngram -lm data/srilm\_interp/lm.\${w}.gz -ppl data/srilm/dev.txt | paste -s **done** | sort -k15,15g > data/srilm\_interp/perplexities.txt # for basic decoding, let's use only a trigram LM [ -d data/lang\_test/ ] && rm -rf data/lang\_test cp -R data/lang data/lang\_test lm=\$(cat data/srilm/perplexities.txt | grep 3gram | head -n1 | awk '{print \$1}'\ local/arpa2G.sh \$1m data/lang\_test data/lang\_test -uu-:---F1 **prepare\_lm.sh** 29% L24 (Shell-script[bash])------

#### local/arpa2G.sh

```
•
                                     Terminal
$decompress | \
 grep -v '<s> <s>' | grep -v '</s> <s>' | grep -v '</s> </s> ' |
 arpa2fst - | \
 fstprint | \
 utils/eps2disambig.pl | \
 utils/s2eps.pl | \
  fstcompile --isymbols=$langdir/words.txt \
  --osymbols=$langdir/words.txt --keep_isymbols=false --keep_osymbols=false | \
 fstrmepsilon | fstarcsort --sort_type=olabel > $destdir/G.fst || exit 1
fstisstochastic $destdir/G.fst || true;
if $cleanup; then
  rm $destdir/lm_tmp.gz 2>/dev/null || true;
fi
exit 0
-uu-:---F1
           arpa2G.sh
                           Bot L108
                                      (Shell-script[bash])---
```

# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building

#### Acoustic model training

- Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring
- Basic DNN system building
- Going beyond the basics



### GMM Training (1)



```
•
                                 Terminal
#!/bin/bash
# Copyright 2012 Johns Hopkins University (Author: Daniel Povey)
# Apache 2.0
# Begin configuration.
stage=-4 # This allows restarting after partway, when something when wrong.
config=
cmd=run.pl
scale_opts="--transition-scale=1.0 --acoustic-scale=0.1 --self-loop-scale=0.1"
realign_iters="10 20 30";
num_iters=35  # Number of iterations of training
max_iter_inc=25 # Last iter to increase #Gauss on.
beam=10
careful=false
retry_beam=40
boost_silence=1.0 # Factor by which to boost silence likelihoods in alignment
power=0.25 # Exponent for number of gaussians according to occurrence counts
-uu-:---F1 align_si.sh
                        Top L1
                                  (Shell-script[bash])-----
-uu-:---F1 train_mono.sh Top L1
                                   (Shell-script[bash])------
Indentation setup for shell type bash
```
# cluster-phones, compile-questions, build-tree

#### • • •

Terminal

#### 

echo "\$0: getting questions for tree-building, via clustering"
# preparing questions, roots file...
cluster-phones \$context\_opts \$dir/treeacc \$lang/phones/sets.int \
 \$dir/questions.int 2> \$dir/log/questions.log || exit 1;
cat \$lang/phones/extra\_questions.int >> \$dir/questions.int
 compile-questions \$context\_opts \$lang/topo \$dir/questions.int \
 \$dir/questions.gst 2>\$dir/log/compile\_questions.log || exit 1;

#### echo "\$0: building the tree"

\$cmd \$dir/log/build\_tree.log \
 build-tree \$context\_opts --verbose=1 --max-leaves=\$numleaves \
 --cluster-thresh=\$cluster\_thresh \$dir/treeacc \$lang/phones/roots.int \
 \$dir/questions.qst \$lang/topo \$dir/tree || exit 1;

\$cmd \$dir/log/init\_model.log \
gmm-init-model --write-occs=\$dir/1.occs \
\$dir/tree \$dir/treeacc \$lang/topo \$dir/1.mdl || exit 1;

-uu-:---F1 **train\_deltas.sh** 48% L94 (Shell-script[bash])------

### GMM Training (4)

• • •	Terminal	
<pre>### Triphone + LDA and M # Training echo "Starting LDA+MLLT steps/align_si.shnj data/train data/lang</pre>	LLT training." \$njcmd "\$train_cmd exp/tri2a exp/tri2a_0	d" \ _ali
steps/train_lda_mllt.sh splice-opts "left-	cmd "\$train_cmd" context=3right-con	∖ ntext=3" ∖
uu-:F1 <b>run.sh</b>	54% L116 (Shell-sc	cript[bash])
<pre>### Triphone + LDA and M # Training echo "Starting SAT+FMLLR share (align ai ab</pre>	LLT + SAT and FMLLR training."	л п. т.
<pre>steps/align_si.shnj use-graphs true da steps/train_sat.shcmd data/train data/lang echo "SAT+FMLLR training</pre>	<pre>\$njcmd "\$train_cmd ta/train data/lang exp "\$train_cmd" 4200 400 exp/tri2b_ali exp/tr done "</pre>	xp/tri2b exp/tri2b_ali 0000 \ ri3b

-uu-:---F1 **run.sh** 67% L138 (Shell-script[bash])------

### GMM Training (5)

Terminal

#### fi

```
if [ $stage -le 8 ]; then
    echo "Starting SGMM training."
    steps/align_fmllr.sh --nj $nj --cmd "$train_cmd" \
        data/train data/lang exp/tri3b exp/tri3b_ali
```

```
steps/train_ubm.sh --cmd "$train_cmd" \
    600 data/train data/lang exp/tri3b_ali exp/ubm5b2
```

```
steps/train_sgmm2.sh --cmd "$train_cmd" \
```

```
5200 12000 data/train data/lang exp/tri3b_ali exp/ubm5b2/final.ubm exp/s\
gmm2_5b2
```

```
echo "SGMM training done."
```



# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building
  - Acoustic model training
  - Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring
- Basic DNN system building
- Going beyond the basics



## Building HCLG (1)

•	Terminal
<b># !</b>	/bin/bash
# 	Copyright 2010-2012 Microsoft Corporation
# #	2012-2013 Johns Hopkins University (Author: Daniel Povey)
Ŧ	Apache 2.0
# #	This script creates a fully expanded decoding graph (HCLG) that represents all the language-model, pronunciation dictionary (lexicon), context-dependenc\
у,	
-u	i-:F1 <b>mkgraph.sh</b> Top L1 (Shell-script[bash])
if	<pre>[[ ! -s \$lang/tmp/LG.fst    \$lang/tmp/LG.fst -ot \$lang/G.fst    \</pre>
	<pre>\$Lang/tmp/LG.fst -ot \$Lang/L_disambig.fst ]]; then</pre>
	<sup>-</sup> sttablecompose <b>\$lang/L_</b> disambig.fst <b>\$lang</b> /G.fst   fstdeterminizestaruse-l\
og	=true   \
	fstminimizeencoded   fstpushspecial   \
	fstarcsortsort_type=ilabel > <b>\$lang</b> /tmp/LG.fst    <b>exit</b> 1;
	<sup>-</sup> stisstochastic <b>\$lang/tmp/LG.fst    echo</b> "[info]: LG not stochastic."
fi	

-uu-:---F1 **mkgraph.sh** 46% L74 (Shell-script[bash])------

## Building HCLG (2)

#### Terminal

```
clg=$lang/tmp/CLG_${N}_${P}.fst
```

```
if [[ ! -s $clg || $clg -ot $lang/tmp/LG.fst ]]; then
  fstcomposecontext --context-size=$N --central-position=$P \
   --read-disambig-syms=$lang/phones/disambig.int \
   --write-disambig-syms=$lang/tmp/disambig_ilabels_${N}_${P}.int \
   $\lang/tmp/ilabels_${N}_${P} < $\lang/tmp/LG.fst |\</pre>
   fstarcsort --sort_type=ilabel > $clg
  fstisstochastic $clg || echo "[info]: CLG not stochastic."
fi
if [[ ! -s $dir/Ha.fst || $dir/Ha.fst -ot $model
   || $dir/Ha.fst -ot $lang/tmp/ilabels_${N}_${P} ]]; then
 if $reverse; then
   make-h-transducer --reverse=true --push_weights=true \
      --disambig-syms-out=$dir/disambig_tid.int \
-uu-:---F1 mkgraph.sh 53% L87 (Shell-script[bash])-----
```

 $\bullet \quad \bigcirc \quad \bigcirc$ 

## Building HCLG (3)

Terminal

#### if [[ ! -s \$dir/HCLGa.fst || \$dir/HCLGa.fst -ot \$dir/Ha.fst || \ \$dir/HCLGa.fst -ot \$clg ]]; then if \$remove\_oov; then [ ! -f \$lang/oov.int ] && \ echo "\$0: --remove-oov option: no file \$lang/oov.int" && exit 1; clg="fstrmsymbols --remove-arcs=true --apply-to-output=true \$lang/oov.int \$\ clg|' fi fsttablecompose \$dir/Ha.fst "\$clg" | fstdeterminizestar --use-log=true \ | fstrmsymbols \$dir/disambig\_tid.int | fstrmepslocal | \ fstminimizeencoded > \$dir/HCLGa.fst || exit 1; fstisstochastic \$dir/HCLGa.fst || echo "HCLGa is not stochastic" fi if [[ ! -s \$dir/HCLG.fst || \$dir/HCLG.fst -ot \$dir/HCLGa.fst ]]; then

add-self-loops --self-loop-scale=\$loopscale --reorder=true \
 \$model < \$dir/HCLGa.fst > \$dir/HCLG.fst || exit 1;

-uu-:---F1 **mkgraph.sh** 70% L112 (Shell-script[bash])-----

## Building HCLG (4)



-uu-:---F1 **mkgraph.sh** 80% L125 (Shell-script[bash])--------

### **Decoding and Lattice Rescoring**

```
# Graph compilation
# Graph compilation
utils/mkgraph.sh data/lang_test exp/sgmm2_5b2 exp/sgmm2_5b2/graph
utils/mkgraph.sh data/lang_big/ exp/sgmm2_5b2 exp/sgmm2_5b2/graph_big
```

•

-uu-:---F1

```
steps/lmrescore_const_arpa.sh --cmd "$decode_cmd" \
    data/lang_test/ data/lang_big/ data/dev \
    exp/sgmm2_5b2/decode_dev exp/sgmm2_5b2/decode_dev.rescored
```

```
steps/decode_sgmm2.sh --nj $dev_nj --cmd "$decode_cmd" \
          --transform-dir exp/tri3b/decode_dev \
          exp/sgmm2_5b2/graph_big data/dev exp/sgmm2_5b2/decode_dev.big
    echo "SGMM decoding done."
) &
```

run.sh 88% L184 (Shell-script[bash])

### steps/decode\_sgmm2.sh

#### Terminal

#### #!/bin/<mark>bash</mark>

•

# Copyright 2012 Johns Hopkins University (Author: Daniel Povey). Apache 2.0.

# This script does decoding with an SGMM system, with speaker vectors. # If the SGMM system was # built on top of fMLLR transforms from a conventional system, you should # provide the --transform-dir option.

# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building
  - Acoustic model training
  - Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring
- Basic DNN system building
- Going beyond the basics



### steps/lmrescore\_const\_arpa.sh

	Terminal
<pre>f ! cmp -s \$oldlang/words.txt \$new echo "\$0: \$oldlang/words.txt and \$ what you are doing."; fi</pre>	<mark>lang/words.txt; then</mark> \$newlang/words.txt differ: make sure you kno∖
oldlmcommand="fstprojectproject_o	output=true \$oldlm  "
uu-:F1 <b>lmrescore_const_arpa.sh</b>	55% L46 (Shell-script[bash])
<pre>f [ \$stage -le 1 ]; then   \$cmd JOB=1:\$nj \$outdir/log/rescore    lattice-lmrescorelm-scale=-1    "ark:gunzip -c \$indir/lat.JOB.ga    lattice-lmrescore-const-arpa</pre>	elm.JOB.log \ .0 \ z " "\$oldlmcommand" ark:- \  \ lm-scale-1 0 \

fi

-uu-:---F1 **lmrescore\_const\_arpa.sh** 70% L57 (Shell-script[bash])-------

ark:- "\$newlm" "ark,t:|gzip -c>\$outdir/lat.JOB.gz" || exit 1;

# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building
  - Acoustic model training
  - Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring

#### Basic DNN system building

• Going beyond the basics



•		Terminal				
fi						
<pre>if [ \$stage -le 7 ]; then     # having a larger number of speakers is helpful for generalization, and to     # handle per-utterance decoding well (iVector starts at zero).     steps/online/nnet2/copy_data_dir.shutts-per-spk-max 2 data/train_hires \         data/train_hires_max2    exit 1     steps/online/nnet2/extract_ivectors_online.shcmd "\$train_cmd"nj 16\</pre>						
data/tra 1 fi	in_hires_max2 exp/nnet3	3/extractor (	exp/nnet3/ivectors_train    <b>exit</b> ∖			
<pre>if [ \$stage -le 8 ]; then     steps/online/nnet2/extract_ivectors_online.shcmd "\$train_cmd"nj 6 \         data/dev_hires exp/nnet3/extractor exp/nnet3/ivectors_dev    exit 1 fi</pre>						
-uu-:F1	<pre>run_ivector_common.sh</pre>	83% L87	(Shell-script[bash])			
-uu-:F1	<pre>run_ivector_common.sh</pre>	69% L75	(Shell-script[bash])			
-uu-:F1	<pre>run_ivector_common.sh</pre>	5% L24	(Shell-script[bash])			

#### steps/nnet3/tdnn/make\_configs.py

 $\bullet \quad \bullet \quad \bullet$ 

Terminal

/Users/sanjeev/Desktop/Conference Travel/SLTU 2016/Tutorial/kaldi/egs/iban/s5 /exp/nnet3/tdnn\_1/configs: total used in directory 88 available 213368700 drwxr-xr-x 13 sanjeev staff 3 12:39 442 May drwxr-xr-x 26 sanjeev 884 May staff 3 13:26 3 12:35 init.config -rw-r--r--1 sanjeev staff 333 May 1 sanjeev staff 1454 May 3 12:35 layer1.config -rw-r--r--1 sanjeev staff 1185 May 3 12:35 layer2.config -rw-r--r--3 12:35 layer3.config -rw-r--r--1 sanjeev staff 1185 May 1 sanjeev 3 12:35 layer4.config staff 1185 May -rw-r--r--1 sanjeev staff 1185 May 3 12:35 layer5.config -rw-r--r--

-rw-r--r-- 1 sanjeev staff 1185 May 3 12:35 layer6.config
-rw-r--r-- 1 sanjeev staff 1122 May 3 12:35 layer7.config
lrwxr-xr-x 1 sanjeev staff 10 May 3 12:39 lda.mat -> ../lda.mat

lrwxr-xr-x 1 sanjeev staff 29 May 3 12:39 presoftmax\_prior\_scale.vec -\
> ../presoftmax\_prior\_scale.vec

-rw-r--r-- 1 sanjeev staff 140 May 3 12:35 vars

-uuu:%%-F1 **configs** All L5 (Dired by name)------Loading dired...done -uu-:---F1 **run\_tdnn.sh** 29% L40 (Shell-script[bash])------

### steps/nnet3/train\_dnn.py

#### • • •

#### Terminal

```
steps/nnet3/train_dnn.py --stage $train_stage \
    --cmd="$decode_cmd" \
    --trainer.optimization.num-jobs-initial 2 \
    --trainer.optimization.num-jobs-final 4 \
    --trainer.num-epochs 4 \setminus
    --trainer.add-layers-period 1 \setminus
    --feat.online-ivector-dir exp/nnet3/ivectors_train
    --feat.cmvn-opts "--norm-means=false --norm-vars=false" \
    --trainer.num-epochs 2 \setminus
    --trainer.optimization.initial-effective-lrate 0.005 \
    --trainer.optimization.final-effective-lrate 0.0005 \
    --trainer.samples-per-iter 120000
    --cleanup.preserve-model-interval 10 \
    --feat-dir data/train_hires \
    --ali-dir exp/nnet3/tri3b_ali_sp \
    --lang data/lang \
    --dir=$dir || exit 1;
                                        (Shell-script[bash])-
-uu-:---F1 run_tdnn.sh
                            51% L60
```

# Building an STT System with Kaldi

- Data preparation
  - Acoustic model training data
  - Pronunciation lexicon
  - Language model training data
- Basic GMM system building
  - Acoustic model training
  - Language model training
- Basic Decoding
  - Creating a static decoding graph
  - Lattice rescoring
- Basic DNN system building
- Going beyond the basics





#### Advanced Methods:

#### Staying Ahead in the STT Game

- STT technology is advancing very rapidly
  - Amazon, Apple, Baidu, Facebook, Google, Microsoft
- Kaldi leads and keeps up with major innovations
  - From SGMMs to DNN (2012)
  - From "English" to low-resource languages (2013)
  - From CPUs to GPUs (2014)
  - From close-talking to far-field microphones (2015)
  - From well-curated to "wild type" corpora (2016)
- A preview of some upcoming developments

#### Advanced Methods: Staying Ahead in the STT Game



- STT technology is advancing very rapidly
  - Amazon, Apple, Baidu, Facebook, Google, Microsoft
- Kaldi leads and keeps up with major innovations

#### - From SGMMs to DNN (2012)

- From "English" to low-resource languages (2013)
- From CPUs to GPUs (2014)
- From close-talking to far-field microphones (2015)
- From well-curated to "wild type" corpora (2016)
- A preview of some upcoming developments

#### **Deep Neural Networks for STT**

INTERSPEECH 2011



#### Conversational Speech Transcription Using Context-Dependent Deep Neural Networks

Frank Seide<sup>1</sup>, Gang Li,<sup>1</sup> and Dong Yu<sup>2</sup>

<sup>1</sup>Microsoft Research Asia, Beijing, P.R.C. <sup>2</sup>Microsoft Research, Redmond, USA {fseide,ganl,dongyu}@microsoft.com

#### Abstract

We apply the recently proposed Context-Dependent Deep-Neural-Network HMMs, or CD-DNN-HMMs, to speech-to-text transcription. For single-pass speaker-independent recognition on the RT03S Fisher portion of phone-call transcription benchmark (Switchboard), the word-error rate is reduced from 27.4%, obtained by discriminatively trained Gaussian-mixture HMMs, to 18.5%—a 33% relative improvement.

CD-DNN-HMMs combine classic artificial-neural-network HMMs with traditional tied-state triphones and deep-beliefnetwork pre-training. They had previously been shown to reduce errors by 16% relatively when trained on tens of hours of data using hundreds of tied states. This paper takes CD-DNN-HMMs further and applies them to transcription using over 300 hours of training data, over 9000 tied states, and up to 9 hidden layers, and demonstrates how sparseness can be exploited. ers), and task (from voice queries to speech-to-text transcription). This is demonstrated on a publicly available benchmark, the Switchboard phone-call transcription task (2000 NIST Hub5 and RT03S sets). We should note here that ANNs have been trained on up to 2000 hours of speech before [7], but with much fewer output units (monophones) and fewer hidden layers.

Second, we advance the CD-DNN-HMMs by introducing weight sparseness and the related learning strategy and demonstrate that this can reduce recognition error or model size.

Third, we present the statistical view of the multi-layer perceptron (MLP) and DBN and provide empirical evidence for understanding which factors contribute most to the accuracy improvements achieved by the CD-DNN-HMMs.

> 2. The Context-Dependent Deep Neural Network HMM

#### **DNN Acoustic Models for the Masses**

- Nontrivial to get the DNN models to work well
  - Design decisions: # layers, # nodes, # outputs, type of nonlinearity, training criterion
  - Training art: learning rates, regularization, update stability (max change), data randomization, # epochs
  - Computational art: matrix libraries, memory mgmt
- Kaldi recipes provide a robust starting point

Corpus	Training Speech	SGMM WER	DNN WER
BABEL Pashto	10 hours	69.2%	67.6%
BABEL Pashto	80 hours	50.2%	42.3%
Fisher English	2000 hours	15.4%	10.3%

### Advanced Methods: Staying Ahead in the STT Game

- STT technology is advancing very rapidly
  - Amazon, Apple, Baidu, Facebook, Google, Microsoft
- Kaldi leads and keeps up with major innovations
   From SGMMs to DNN (2012)
  - From "English" to low-resource languages (2013)
  - From CPUs to GPUs (2014)
  - From close-talking to far-field microphones (2015)
  - From well-curated to "wild type" corpora (2016)
- A preview of some upcoming developments

### Low-Resource STT for the Masses

- Kaldi provides language-independent recipes
  - Typical BABEL Full LP condition
    - 80 hours of transcribed speech, 800K words of LM text, 20K word pronunciation lexicon
  - Typical BABEL Limited LP condition
    - 10 hours of transcribed speech, 100K words of LM text, 6K word pronunciation lexicon

Language	Cantonese		Tagalog		Pashto		Turkish	
Speech	80h	10h	80h	10h	80h	10h	80h	10h
CER/WER	48.5%	61.2%	46.3%	61.9%	50.7%	63.0%	51.3%	65.3%
ATWV	0.47	0.26	0.56	0.28	0.46	0.25	0.52	0.25

### Advanced Methods: Staying Ahead in the STT Game

- STT technology is advancing very rapidly
  - Amazon, Apple, Baidu, Facebook, Google, Microsoft
- Kaldi leads and keeps up with major innovations
  - From SGMMs to DNN (2012)
  - From "English" to low-resource languages (2013)
  - From CPUs to GPUs (2014)
  - From close-talking to far-field microphones (2015)
  - From well-curated to "wild type" corpora (2016)
- A preview of some upcoming developments

### Parallel (GPU-based) Training

- Original neural network training algorithms were inherently sequential (e.g. SGD)
- Scaling up to "big data" becomes a challenge
- Several solutions have emerged recently
  - 2009: Delayed SGD (Yahoo!)
  - 2011: Lock-free SGD (Hogwild! U Wisconsin)
  - 2012: Gradient averaging (DistBelief, Google)
  - 2014: Model averaging (NG-SGD, Kaldi)

#### PARALLEL TRAINING OF DNNS WITH NATURAL GRA-DIENT AND PARAMETER AVERAGING

#### Daniel Povey, Xiaohui Zhang & Sanjeev Khudanpur

Center for Language and Speech Processing & Human Language Technology Center of Excellence, The Johns Hopkins University, Baltimore, MD 21218, USA {dpovey@gmail.com}, {xiaohui,khudanpur@jhu.edu}

#### ABSTRACT

We describe the neural-network training framework used in the Kaldi speech recognition toolkit, which is geared towards training DNNs with large amounts of training data using multiple GPU-equipped or multi-core machines. In order to be as hardware-agnostic as possible, we needed a way to use multiple machines without generating excessive network traffic. Our method is to average the neural network parameters periodically (typically every minute or two), and redistribute the averaged parameters to the machines for further training. Each machine sees different data. By itself, this method does not work very well. However, we have another method, an approximate and efficient implementation of Natural Gradient for Stochastic Gradient Descent (NG-SGD), which seems to allow our periodicaveraging method to work well, as well as substantially improving the convergence of SGD on a single machine.

### Model Averaging with NG-SGD

- Train DNNs with large amount of data
  - Utilize a cluster of CPUs or GPUs
  - Minimize network traffic (esp. for CPUs)
- Solution: exploit data parallelization
  - Update model in parallel over many mini-batches
  - Infrequently average models (parameters)
- Use "Natural-Gradient" SGD for model updating
  - Approximates conditioning via inverse Fisher matrix
  - Improves convergence even without parallelization

#### **Parallelization Matters!**



- Typically, a GPU is 10x faster than a 16 core CPU
- Linear speed-up till ca 4 GPUs, then diminishing

#### Advanced Methods: Staying Ahead in the STT Game

- STT technology is advancing very rapidly
  - Amazon, Apple, Baidu, Facebook, Google, Microsoft
- Kaldi leads and keeps up with major innovations
  - From SGMMs to DNN (2012)
  - From "English" to low-resource languages (2013)
  - From CPUs to GPUs (2014)
  - From close-talking to far-field microphones (2015)
  - From well-curated to "wild type" corpora (2016)
- A preview of some upcoming developments

### IARPA's Open Challenge

 Automatic speech recognition software that works in a variety of acoustic environments and recording scenarios is a holy grail of the speech research community.



 IARPA's Automatic Speech recognition In Reverberant Environments (ASpIRE) Challenge is seeking that grail.

### Rules of the ASpIRE Challenge

- 15 hours of speech data were posted on the IARPA website
  - Multi-microphone recordings of conversational English
  - 5h development set (dev), 10h development-test set (dev-test)
  - Transcriptions provided for dev, only scoring for dev-test output
  - For training data selection, system development and tuning
- 12 hours of new speech data during the evaluation period
  - Far-field speech (eval) from noisy, reverberant rooms
  - Single-microphone or multi-microphone conditions
- Word error rate is the measure of performance
  - Single-microphone submissions were due on 02/18/2015
  - Results were officially announced on 09/10/2015

### Examples of ASpIRE Audio

- Typical sample
  - Suggested by Dr. Mary Harper
- Almost manageable
  - Easy for humans, 26% errors for ASR
  - Somewhat hard
    - Easy for humans, 41% errors for ASR
    - Much harder
      - Not easy for humans, 60% errors for ASR
    - \*#@!!#% #%^^
      - Very hard for humans, no ASR output



#### Kaldi ASR Improvements for ASpIRE

- Time delay neural networks (TDNN)
  - A way to deal with long acoustic-phonetic context
  - A structured alternative to deep/recurrent neural nets
- Data augmentation with simulated reverberations
  - A way to mitigate channel distortions not seen in training
  - A form of multi-condition training of ASR models
- i-vector based speaker & environment adaptation
  - A way to deal with speaker & channel variability
  - Adapted [with a twist] from Speaker ID systems

#### Kaldi ASR Improvements, ASpIRE++

- Pronunciation and inter-word silence modeling
  - Inspired by pronunciation-prosody interactions
  - A simple context-dependent model of inter-word silence
- Recurrent neural network language models (RNNLM)
  - A (known) way to model long-range word dependencies
  - Incorporated post-submission into JHU ASpIRE system
- Ongoing Kaldi investigations that hold promise
  - Semi-supervised discriminative training of (T)DNNs
  - Long short-term memory (LSTM) acoustic models
  - Connectionist temporal classification (CTC) models

#### **Time Delay Neural Networks**

(See our paper at INTERSPEECH 2015 for details)

#### A time delay neural network architecture for efficient modeling of long temporal contexts

Vijayaditya Peddinti<sup>1</sup>, Daniel Povey<sup>1,2</sup>, Sanjeev Khudanpur<sup>1,2</sup>

<sup>1</sup>Center for Language and Speech Processing & <sup>2</sup>Human Language Technology Center of Excellence Johns Hopkins University, Baltimore, MD 21218, USA

vijay.p,khudanpur@jhu.edu, dpovey@gmail.com

#### Abstract

Recurrent neural network architectures have been shown to efficiently model long term temporal dependencies between acoustic events. However the training time of recurrent networks is higher than feedforward networks due to the sequential nature of the learning algorithm. In this paper we propose a time delay neural network architecture which models long term temporal dependencies with training times comparable to standard feed-forward DNNs. The network uses sub-sampling to reduce computation during training. On the Switchboard task we show a relative improvement of 6% over the baseline DNN model. We present results on several LVCSR tasks with training data ranging from 3 to 1800 hours to show the effectiveness of the TDNN architecture in learning wider temporal dependencies in both small and large data scenarios.

**Index Terms**: time delay neural networks, acoustic modeling, recurrent neural networks

be reduced, while ensuring that information from all time steps in the input context is processed by the network.

Neural network architectures have been shown to benefit from speaker adaptation. However, speaker adaptation techniques like fMLLR [4] require two passes of decoding. The 2pass decoding strategy is difficult to use in online speech recognition applications. iVectors which capture both speaker and environment specific information have been shown to be useful for instantaneous and discriminative adaptation of the neural network [5, 6]. In this paper we use iVector based neural network adaptation.

The paper is organized as follows. Section 2 mentions relevant work, Section 3 describes the neural network architecture and training recipe in greater detail. Section 4 describes the experimental setup. Section 5 presents and analyzes the results primarily on the Switchboard [7] task. It also presents results on other LVCSR tasks which have 3-1800 hours of training data. Section 6 presents the conclusions and the future work.

#### A 28 Year Old Idea, Resurrected



Alex Waibel, Kevin Lang, et al (1987)
# Improved ASR on Several Data Sets

Standard ASR Test Sets	Size	DNN	TDNN	Rel. Δ
Wall Street Journal	80 hrs	6.6%	6.2%	5%
TED-LIUM	118 hrs	19.3%	17.9%	7%
Switchboard	300 hrs	15.5%	14.0%	10%
Libri Speech	960 hrs	5.2%	4.8%	7%
Fisher English	1800 hrs	22.2%	21.0%	5%

- Consistent 5-10% reduction in word error rate (**WER**) over DNNs on most datasets, including conversational speech.
- TDNN training speeds are on par with DNN, and nearly an order of magnitude faster than RNN

ASpIRE (Fisher Training) 1800 hrs 47.7% 47.6%

### Data Augmentation for ASR Training

(See our paper at INTERSPEECH 2015 for details)

### Audio Augmentation for Speech Recognition

Tom Ko<sup>1</sup>, Vijayaditya Peddinti<sup>2</sup>, Daniel Povey<sup>2,3</sup>, Sanjeev Khudanpur<sup>2,3</sup>

<sup>1</sup>Huawei Noah's Ark Research Lab, Hong Kong, China <sup>2</sup>Center for Language and Speech Processing & <sup>3</sup>Human Language Technology Center of Excellence, Johns Hopkins University, Baltimore, MD, 21218, USA {tomkocse,dpovey}@gmail.com, {vijay.p,khudanpur}@jhu.edu

#### Abstract

Data augmentation is a common strategy adopted to increase the quantity of training data, avoid overfitting and improve robustness of the models. In this paper, we investigate audio-level speech augmentation methods which directly process the raw signal. The method we particularly recommend is to change the speed of the audio signal, producing 3 versions of the original signal with speed factors of 0.9, 1.0 and 1.1. The proposed technique has a low implementation cost, making it easy to adopt. We present results on 4 different LVCSR tasks with training data ranging from 100 hours to 1000 hours, to examine the effectiveness of audio augmentation in a variety of data scenarios. An average relative improvement of 4.3% was observed across the 4 tasks.

Index Terms: speech recognition, data augmentation, deep neural network

#### 2. Audio perturbation

In this section we describe a speed-perturbation technique for data augmentation and compare it with the existing augmentation technique VTLP [3]. Speed perturbation produces a warped time signal. Given an audio signal x(t), time warping by a factor  $\alpha$  gives the signal  $x(\alpha t)$ . It can be seen from the Fourier transform of  $x(\alpha t)$ ,  $\alpha^{-1}\hat{x}(\alpha^{-1}\omega)$ , that the warping factor produces shifts in the frequency components of the  $\hat{x}(\omega)$ by an amount proportional to frequency  $\omega$ . In [8] it was shown that this corresponds approximately to a shift of the spectrum in the mel spectrogram, since the mel scale is approximately logarithmic. It can be seen that these changes in the mel spectrogram are similar to those produced using VTLP. However, unlike VTLP, speed perturbation results in a change in the duration of the signal which also affects the number of frames in the utterance.

# Simulating Reverberant Speech for Multi-condition (T)DNN Training

- Simulate ca 5500 hours of reverberant, noisy data from 1800 hours of the Fisher English CTS corpus
  - Replicate each of the ca 21,000 conversation sides 3 times
  - Randomly change the sampling rate [up to ±10%]
  - Convolve each conversation side with one of 320 real-life room impulse responses (RIR) chosen at random
  - Add noise to the signal (when available with the RIR)
- Generate (T)DNN training labels from clean speech
   Align "pre-reverb" speech to ca 7500 CD-HMM states
- Train DNN and TDNN acoustic models

Cross-entropy training followed by sequence training

# **Result of Data Augmentation**

Acoustic Model	Data Augmentation	Dev WER
TDNN A (230 ms)	None (1800h, clean speech)	47.6%
TDNN A (230 ms)	+ 3 x (reverberation + noise)	31.7%
TDNN B (290 ms)	+ 3 x (reverberation + noise)	30.8%
TDNN A (230 ms)	+ sampling rate perturbation	31.0%
TDNN B (290 ms)	+ sampling rate perturbation	31.1%

- Data augmentation with simulated reverberation is beneficial
  - Likely to be a very important reason for relatively good performance
- Sampling rate perturbation didn't help much on ASpIRE data
- Sequence training helped reduce WER on the dev set
  - Required modifying the sMBR training criterion to realize gains
  - But the gains did not carry over to dev-test set

### i-vectors for Speaker Compensation

(See our paper at INTERSPEECH 2015 for details)

### Reverberation robust acoustic modeling using i-vectors with time delay neural networks

*Vijayaditya Peddinti*<sup>1</sup>, *Guoguo Chen*<sup>1</sup>, *Daniel Povey*<sup>1,2</sup>, *Sanjeev Khudanpur*<sup>1,2</sup>

<sup>1</sup>Center for language and speech processing & <sup>2</sup>Human Language Technology Center of Excellence The Johns Hopkins University, Baltimore, MD 21218, USA

{vijay.p,guoguo,khudanpur}@jhu.edu, dpovey@gmail.com

#### Abstract

In reverberant environments there are long term interactions between speech and corrupting sources. In this paper a time delay neural network (TDNN) architecture, capable of learning long term temporal relationships and translation invariant representations, is used for reverberation robust acoustic modeling. Further, iVectors are used as an input to the neural network to perform instantaneous speaker and environment adaptation, providing 10% relative improvement in word error rate. By subsampling the outputs at TDNN layers across time steps, training time is reduced. Using a parallel training algorithm we show that the TDNN can be trained on  $\sim 5500$  hours of speech data in 3 days using up to 32 GPUs. The TDNN is shown to provide results competitive with state of the art systems in the IARPA ASpIRE challenge, with 27.7% WER on the  $dev\_test$  set. Index Terms: far field speech recognition, time delay neural networks, reverberation

such a wide temporal context enables the network to deal with late reverberations.

iVectors which capture both speaker and environment specific information have been shown to be useful for rapid adaptation of the neural network [6, 7, 8]. iVector based adaptation has also been shown to be effective in reverberant environments [9]. In this paper we use this adaptation technique.

We show experimental results on the ASpIRE far-field speech recognition challenge held by IARPA [10]. This challenge uses the English portion of the Fisher database [11] for acoustic and language model training. We show that in this large data scenario the proposed network architecture, combined with a parallel training technique [12], can train on multi-condition training data of ~ 5500 hours, using up to 32 GPUs, in 3 days.

Using the TDNN architecture helps us to achieve results close to those of the best combined system submitted to the AS-pIRE challenge, while using only a single system. Our system was able to achieve 27.7% WER on the *dev-test* set, while the best system achieved 27.2% WER.

1 Introduction

### Using i-vectors Instead of fMLLR and using unnormalized MFCCs to compute i-vectors

- 100-dim i-vectors are appended to MFCC inputs of the TDNN
  - i-vectors are computed from raw MFCCs (i.e. no mean subtraction etc)
  - UBM posteriors however use MFCCs normalized over a 6 sec window
- i-vectors are computed for each training utterance
  - Increases speaker- and channel variability seen in training data
  - May model transient distortions? e.g. moving speakers, passing cars
- i-vectors are calculated for every ca 60 sec of test audio
  - UBM prior is weighted 10:1 to prevent overcompensation
  - Weight of test statistics is capped at 75:1 relative to UBM statistics

Speaker Compensation Method	Dev WER
TDNN without i-vectors	34.8%
+ i-vectors (from all frames)	33.8%
+ i-vectors (from reliable speech frames)	30.8%

### **Pronunciation and Silence Probabilities**

(See our paper at INTERSPEECH 2015 for details)

### **PRONUNCIATION AND SILENCE PROBABILITY MODELING FOR ASR**

*Guoguo Chen*<sup>1</sup>, *Hainan Xu*<sup>1</sup>, *Minhua Wu*<sup>1</sup>, *Daniel Povey*<sup>1,2</sup>, *Sanjeev Khudanpur*<sup>1,2</sup>

<sup>1</sup>Center for Language and Speech Processing <sup>2</sup>Human Language Technology Center of Excellence The Johns Hopkins University, Baltimore, MD 21218, USA

guoguo@jhu.edu, hxu31@jhu.edu, mwu56@jhu.edu, dpovey@gmail.edu, khudanpur@jhu.edu

### Abstract

In this paper we evaluate the WER improvement from modeling pronunciation probabilities and word-specific silence probabilities in speech recognition. We do this in the context of Finite State Transducer (FST)-based decoding, where pronunciation and silence probabilities are encoded in the lexicon (L) transducer. We describe a novel way to model word-dependent silence probabilities, where in addition to modeling the probability of silence following each individual word, we also model the probabilities are estimated from aligned training data, with suitable smoothing. We conduct our experiments on four commonly used automatic speech recognition datasets, namely Wall Street Journal, Switchboard, TED-LIUM, and Librispeech. The Implicit pronunciation modeling relies on the underlying acoustic-phonetic models to account for pronunciation variations, and therefore removes the necessity to explicitly determine and represent them in the lexicon. In some methods, acoustic model parameters of a phoneme (e.g., Gaussian densities) are tied with those of its alternative realizations, thus capturing alternative pronunciations [3, 12, 13]. Others view pronunciations as a bundle of features, and pronunciation variation is viewed as feature-change or asynchrony [14, 15].

While variability in the pronunciation of individual words has been studied extensively, relatively little has been studied about inter-word silence and its dependence on the prosodic and syntactic structure of the utterance. In [3, 10], for instance, three types of silence are permitted following each pronunciation in the lexicon: a zero-silence, a short pause and a long silence. It is

### **Trigram-like Inter-word Silence Model**



# Is "Prosody" Finally Helping STT?

Task	Test Set	Baseline	+ Sil/Pron
WSJ	Eval 92	4.1	3.9
Switchboard	Eval 2000	20.5	20.0
TED-LIUM	Test	18.1	17.9
Libri Crossh	Test Clean	6.6	6.6
Libri Speech	Test Other	22.9	22.5

• Modeling pronunciation and silence probabilities yields modest but consistent improvement on many large vocabulary ASR tasks

Pronunciation/Silence Probabilities	Dev WER
No probabilities in the lexicon	32.1%
+ pronunciation probabilities	31.6%
+ inter-word silence probabilities	30.8%

## Recurrent Neural Network based Language Models

(See our paper at INTERSPEECH 2010 for the first "convincing" results)

### **Recurrent neural network based language model**

Tomáš Mikolov<sup>1,2</sup>, Martin Karafiát<sup>1</sup>, Lukáš Burget<sup>1</sup>, Jan "Honza" Černocký<sup>1</sup>, Sanjeev Khudanpur<sup>2</sup>

<sup>1</sup>Speech@FIT, Brno University of Technology, Czech Republic
<sup>2</sup> Department of Electrical and Computer Engineering, Johns Hopkins University, USA {imikolov, karafiat, burget, cernocky}@fit.vutbr.cz, khudanpur@jhu.edu

### Abstract

A new recurrent neural network based language model (RNN LM) with applications to speech recognition is presented. Results indicate that it is possible to obtain around 50% reduction of perplexity by using mixture of several RNN LMs, compared to a state of the art backoff language model. Speech recognition experiments show around 18% reduction of word error rate on the Wall Street Journal task when comparing models trained on the same amount of data, and around 5% on the much harder NIST RT05 task, even when the backoff model is trained on much more data than the RNN LM. We provide ample empirical evidence to suggest that connectionist language models are superior to standard n-gram techniques, except their high computational (training) complexity.



# RNN LM on ASpIRE Data

Language Model and Rescoring Method	Dev WER
4-gram LM and lattice rescoring	30.8%
RNN-LM and 100-best rescoring	30.2%
RNN-LM and 1000-best rescoring	29.9%
RNN-LM (4-gram approximation) lattice rescoring	29.9%
RNN-LM (6-gram approximation) lattice rescoring	<b>29.8</b> %

- An RNN LM consistently outperforms the N-gram LM
- The Kaldi lattice rescoring appears to cause no loss in performance
  - Approximation entails not "expanding" the lattice to represent each unique history separately
  - When two paths merge in an N-gram lattice, only one s(t) is chosen at random and propagated forward

## The IARPA ASpIRE Leader Board

INI	NOCI	ENTIVE®			1-85	5-CROWDNOV	V · Contact Us · Blog   aspire_iarpa ·	Log Out			
M	ly IC	Products/Services	For Solvers	Challenge Cen	ter Res	ources	About Us Challenge Search	>>>			
AS Rev	pIRE verb	– IARPA Au erant Enviror Details Test Solution	Itomatic S Inments Ch My Solution	peech reg nallenge eaderboard Mee	cogniti	on in					
	B	ASpIRE – IARPA Reverberant Envi AWARD: See details   DEA Source: InnoCentive	Automatic Spe ronments Chal DLINE: 2/18/15   ACTIV hallenge ID: 9933624	ech recognitio lenge E SOLVERS: 160   PO Type: RTP	IN IN ISTED: 11/17/14						
				Share 3	■ 6 Messages	Agreement	t				
9	Solver S	olution Scores Rank U	lser Name	Score							
	_	1 vi	ijaypeddinti	72.20		Rank	Participant	D	ev WER	Syster	n Type
		3	SriramG	70.40							
		4	rhsiao	69.40		1	tsakilidis		27.2%	Comb	ination
		5	burget	68.10 67.60		2					
		7 a	spire_iarpa	65.10		2	rhsiao		27.5%	Comb	ination
		8	falavi	60.00		С	vijavnoddinti		סד דר/	Single	Suctom
		9	vmitra	56.60 46.00		5	vijaypeuulitti		21.1/0	Single	System
	The top 10	submissions made via the Te	st Solution tab appear i	n the leaderboard abo	ve						

http://www.dni.gov/index.php/newsroom/press-releases/210-press-releases-2015/1252-iarpa-announces-winners-of-its-aspire-challenge



### IARPA Announces Winners of its ASpIRE Challenge

### NEWS RELEASE

### FOR IMMEDIATE RELEASE ODNI News Release No. 14-15 September 10, 2015

### IARPA Announces Winners of its ASpIRE Challenge

WASHINGTON – The Intelligence Advanced Research Projects Activity (IARPA), within the Office of the Director of National Intelligence (ODNI), today announced the winners of its speech recognition challenge, *Automatic Speech recognition in Reverberant Environments* (ASpIRE). The winning teams from the Johns Hopkins University, Raytheon BBN Technologies, the Institute for Infocomm Research, and Brno University of Technology will share \$110,000 in prizes. http://www.dni.gov/index.php/newsroom/press-releases/210-press-releases-2015/1252-iarpa-announces-winners-of-its-aspire-challenge



### IARPA Announces Winners of its ASpIRE Challenge

The Multiple Microphone Condition tested accuracy of speech recognition on recordings from six different microphones recording at once.

All of the ASpIRE challenge winners delivered systems with more than a 50% reduction in word error rate (WER) compared to the IARPA baseline system. WER is the standard measure of accuracy for speech recognition systems; lower WER scores indicate more accurate systems.

The winners in the Single Microphone category are:

 the team from the Center for Language and Speech Processing, Johns Hopkins University (Vijayaditya Peddinti, Guoguo Chen, Dr. Daniel Povey, Dr. Sanjeev Khudanpur);

# Performance on Evaluation Data

Participant	Test WER	System Type
Kaldi	44.3%	Single System
BBN (and others)	44.3%	Combination
I <sup>2</sup> R (Singapore)	44.8%	Combination

Acoustic Model	Language Model	Dev WER	Test WER	Eval WER
TDNN B (CE training)	4-gram	30.8%	27.7%	44.3%
TDNN B (sMBR training)	4-gram	29.1%	28.9%	<b>43.9</b> %
TDNN B (CE training)	RNN	29.8%	26.5%	43.4%
TDNN B (sMBR training)	RNN	28.3%	28.2%	43.4%

### Keys to Good Performance on ASpIRE

Time delay neural networks (TDNN)

Deal well with long reverberation times

- i-vector based adaptation compensation
   Deals with speaker & channel variability
- Data augmentation with simulated reverberations
   Deals with channel distortions not seen in training
- Pronunciation and inter-word silence probabilities

- Helpful in adverse acoustic conditions

## The JHU ASpIRE System

### (See our ASRU 2015 paper for details)

### JHU ASPIRE SYSTEM : ROBUST LVCSR WITH TDNNS, I-VECTOR ADAPTATION AND RNN-LMS

Vijayaditya Peddinti<sup>1</sup>, Guoguo Chen<sup>1</sup>, Vimal Manohar<sup>1</sup>, Tom Ko<sup>3</sup> Daniel Povey<sup>1,2</sup>, Sanjeev Khudanpur<sup>1,2</sup>

<sup>1</sup>Center for language and speech processing & <sup>2</sup>Human Language Technology Center of Excellence The Johns Hopkins University, Baltimore, MD 21218, USA <sup>3</sup>Huawei Noah's Ark Research Lab, Hong Kong, China

{vijay.p,guoguo,khudanpur}@jhu.edu, {vimal.manohar91, tomkocse, dpovey}@gmail.com

#### ABSTRACT

Multi-style training, using data which emulates a variety of possible test scenarios, is a popular approach towards robust acoustic modeling. However acoustic models capable of exploiting large amounts of training data in a comparatively short amount of training time are essential. In this paper we tackle the problem of reverberant speech recognition using 5500 hours of simulated reverberant data. We use time-delay neural network (TDNN) architecture, which is capable of tackling long-term interactions between speech and corrupting sources in reverberant environments. By sub-sampling the outputs at TDNN layers across time steps, training time is substantially reduced. Combining this with distributed-optimization we show that the TDNN can be trained in 3 days using up to 32 GPUs. Further, iVectors are used as an input to the neural network to perform iniVectors which capture both speaker and environment specific information have been shown to be useful for rapid adaptation of the neural network [4, 5, 6]. iVector based adaptation has also been shown to be effective in reverberant environments [7]. In this paper we use this adaptation technique.

We show experimental results on the ASpIRE far-field speech recognition challenge held by IARPA [8]. This challenge uses the English portion of the Fisher database [9] for acoustic and language model training. We show that in this large data scenario the proposed network architecture, combined with a distributed optimization technique [10], can train on multi-condition training data of ~ 5500 hours, using up to 32 GPUs, in 3 days.

Using the TDNN architecture helps us to achieve results close to those of the best combined system submitted to the ASpIRE chal-

### Semi-supervised MMI Training

(See our paper at INTERSPEECH 2015 for details)

### Semi-supervised Maximum Mutual Information Training of Deep Neural Network Acoustic Models

Vimal Manohar\*, Daniel Povey\*<sup>†</sup>, Sanjeev Khudanpur\*<sup>†</sup>

\*Center for Language and Speech Processing <sup>†</sup> Human Language Technology Center of Excellence The Johns Hopkins University, Baltimore, MD 21218, USA

vmanoha1@jhu.edu, danielpovey@gmail.com, khudanpur@jhu.edu

#### Abstract

Maximum Mutual Information (MMI) is a popular discriminative criterion that has been used in supervised training of acoustic models for automatic speech recognition. However, standard discriminative training is very sensitive to the accuracy of the transcription and hence its implementation in a semisupervised setting requires extensive filtering of data. We will show that if the supervision transcripts are not known, the natural analogue of MMI is to minimize the conditional entropy of the lattice of possible transcripts of the data. This is equivalent to the weighted average of MMI criterion over different reference transcripts, taking those reference transcripts and their weighting from the lattice itself. In this paper we describe experiments where we applied this method to the semi-supervised training of Deep Neural Network acoustic models. In our experimental setup, the proposed method gives up to 0.5% absolute untranscribed data for cross-entropy training using a slightly different method which we will described here. We avoid the untranscribed data "polluting" the last layer of the network by giving it a separate final layer, using ideas inspired by multilingual DNN training [15, 16].

Discriminative training is very sensitive to the accuracy of the transcripts [17, 18, 19]. Therefore sequence-discriminative self-training methods do not work well without some form of confidence-based filtering, as used in [20, 21, 17, 22]. However, we show in this paper that by using an alternative objective function, Negative Conditional Entropy (NCE) on the untranscribed portion of the data, we can obtain improvements from untranscribed data without filtering. Entropy minimization has previously been used as an objective for semi-supervised learning in a facial recognition problem [23] and for sequencediscriminative training of GMM acoustic models for speech



### Semi-Supervised Sequence Training

- Sequence training improves substantially over basic cross-entropy training of DNN acoustic models
- Semi-supervised cross-entropy training by adding unlabeled data – also improves substantially over basic cross-entropy training on labeled data
- But semi-supervised sequence training is "tricky"
  - Sensitivity to incorrect transcription seems greater
  - Confidence-based filtering or weighting must be applied
  - Empirical results are not very satisfactory

### Semi-supervised Sequence Training: without committing to a single transcription

View MMI training as minimizing a conditional entropy

$$I(W \land O; \theta) = \frac{1}{T} \sum_{t=1}^{T} \log \frac{P(O_t | W_t; \theta)}{P(O_t; \theta)} = \frac{1}{T} \sum_{t=1}^{T} \log \frac{P(O_t | W_t; \theta)}{\sum_{W'} P(O_t | W'; \theta) P(W')}$$

$$I(W \land O; \theta) = H(W) - H(W|O; \theta) = H(W) - \frac{1}{T} \sum_{t=1}^{T} H(W|O_t; \theta)$$

- The latter does not require committing to a single  $W_t$ 
  - Well suited for unlabeled speech
  - Entails computing a sum over all W's in the lattice

## Computing Lattice Entropy Using Expectation Semi-rings

• How to efficiently compute  $-H(W|O_t;\theta) = \sum_{\pi \in L} P(\pi) \log P(\pi)$ 



- Take inspiration from the computation of  $Z(O_t; \theta) = \sum P(\pi)$
- Replace arc-probabilities  $p_i$  with the pair ( $p_i$ ,  $p_i \log\{p_i\}$ )

Semi-ring Element & Operators	( p , p×log{p} )
$(p_1, p_1 \log\{p_1\}) + (p_2, p_2 \log\{p_2\})$	$(p_1+p_2, p_1\log\{p_1\}+p_2\log\{p_2\})$
$(p_1, p_1 \log\{p_1\}) \times (p_2, p_2 \log\{p_2\})$	$(p_1p_2, p_1p_2\log\{p_2\}+p_2p_1\log\{p_1\})$

Semi-supervised Sequence Training: Key Details Needed to Make it Work

- View training criterion as MCE instead of MMI
  - i.e. arg min  $H(W|O;\vartheta)$  instead of arg max  $I(W \land O;\vartheta)$
  - Efficiently compute H(W/O;ϑ) for the lattice, and its gradient, via Baum Welch with special semi-rings
- Use separate output (soft-max) layers in the DNN for labeled and unlabeled data

Inspired by multilingual DNN training methods

• Use a slightly different "prior" for converting DNN posterior probabilities to acoustic likelihoods

# Results for Semi-Supervised MMI on Fisher English CTS

ata	DNN Training Method (hours of speech)	Dev WER	Test WER
ed d	Cross-Entropy Training (100h labeled)	32.0	31.2
bele	CE (100h labeled + 250h self-labeled)	30.6	29.8
unla	CE (100h labeled + 250h weighted)	30.5	29.8

**Better use** of unlabeled data

Known use of

Sequence Training (100h labeled)	29.6	28.5
Seq Training (100h labeled +250h weighted)	29.9	28.8
Seq Training (100h labeled + 250h MCE)	29.4	28.1
Sequence Training (350h labeled)	28.5	27.5

- Recovers about 40% of the supervised training gain
  - Investigation underway for 2000h of unlabeled speech
- Repeatable results on BABEL datasets with 10h supervised training + 50-70h unsupervised

## Advanced Methods: Staying Ahead in the STT Game



- STT technology is advancing very rapidly
  - Amazon, Apple, Baidu, Facebook, Google, Microsoft
- Kaldi leads and keeps up with major innovations
  - From SGMMs to DNN (2012)
  - From "English" to low-resource languages (2013)
  - From CPUs to GPUs (2014)
  - From close-talking to far-field microphones (2015)
  - From well-curated to "wild type" corpora (2016)
- A preview of some upcoming developments

# Heterogeneous Training Corpora

- Transcribed speech from different collections are not easy to merge for STT training
  - Genre and speaking style differences
  - Different channel conditions
  - Slightly different transcription conventions
- Typical result: the corpus matched to test data gives best STT results; others don't help, sometimes hurt!
- SCALE 2015 case study with Pashto CTS
  - Collected in country, and transcribed, by same vendor
  - Roughly 80 hours each in the
    - Appen LILA corpus and IARPA BABEL corpus
  - Pronunciation lexicon to cover transcripts; same phone set

### A Study in Pashto

(A manuscript is in preparation for future publication)

### USING OF HETEROGENEOUS CORPORA FOR TRAINING OF AN ASR SYSTEM

Jan Trmal<sup>1,2</sup>, Gaurav Kumar<sup>1,2</sup>, Vimal Manohar<sup>1</sup>, Sanjeev Khudanpur<sup>1,2</sup>, Matt Post<sup>2</sup>, Paul McNamee<sup>2</sup>

<sup>1</sup>Center for Language and Speech Processing <sup>2</sup>Human Language Technology Center of Excellence The Johns Hopkins University, Baltimore, MD 21218, USA

#### ABSTRACT

The paper summarizes the development of the LVCSR system built as a part of the Pashto speech-translation system at the SCALE (Summer Camp for Applied Language Exploration) 2015 workshop on "Speech-to-text-translation for low-resource languages". The Pashto language was chosen as a good "proxy" low-resource language, exhibiting multiple phenomena which make the speech-recognition and and speech-to-text-translation systems development hard.

Even when the amount of data is seemingly sufficient, given the fact that the data originates from multiple sources, the preliminary experiments reveal that there is little to no benefit in merging (concatenating) the corpora and more elaborate ways of making use of all of the data must be worked out.

This paper concentrates only on the LVCSR part and presents a range of different techniques that were found to be useful in order to benefit from multiple different corpora Fig. 1. ARABIC LETTER FARSI YEH and ARABIC LETTER ALEF MAKSURA

layout although the Arabic and Urdu layouts are used as well. The alternative layouts have a majority of the Pashto characters (and people freely substitute those which are missing with visually similar characters). Also, different fonts can have small deficiencies in rendering of glyphs, especially during kerning or joining of characters and the users often try to fix this by substituting a different character that looks better (i.e. closer to the expected shape) in the given context.

A direct impact of this is that a word as a sequence of glyphs (visual representations of characters) can be represented as multiple sequences of unicode codepoints (numerical codes

# A Study in Pashto

- Transcriptions require extensive cross-corpus normalization
- Even after that, language models don't benefit much from corpus pooling
- Simple corpus pooling doesn't improve acoustic modeling either
- DNNs with shared "inner" layers and corpus-specific input and output layers work best

Training Data	Single Model	Interp LM A	olation W LM B	/eights LM T	Interpolated Model
Text A	99.2	0.8	0.2	0.0	92.9
Text B	141.9	0.1	0.8	0.1	140.0
Text T	86.7	0.0	0.0	1.0	86.7

DNN Training Data	STT Word Error Rates		
	Test A	Test B	Test T
Single corpus (matched)	55.4%	<b>46.8%</b>	24.8%
Two corpora (Pashto A + B)	51.9%	48.2%	52.6%

Multi-corpus (A+B) Training	STT Word Error Rates			
Strategy	Test A	Test B	Test T	
Shared DNN layers (except 1)	53.2%	47.4%	27.0%	
Shared DNN layers (except 2)	51.2%	45.0%	25.4%	
+ Optimized Language Model	50.8%	44.8%	25.4%	
+ Duration Modeling	50.4%	44.3%	24.8%	

## Advanced Methods: Staying Ahead in the STT Game



- STT technology is advancing very rapidly
  - Amazon, Apple, Baidu, Facebook, Google, Microsoft
- Kaldi leads and keeps up with major innovations
  - From SGMMs to DNN (2012)
  - From "English" to low-resource languages (2013)
  - From CPUs to GPUs (2014)
  - From close-talking to far-field microphones (2015)
  - From well-curated to "wild type" corpora (2016)
- A preview of some upcoming developments

# **Other Additions and Innovations**

- Semi-supervised (MMI) training
  - Using unlabeled speech to augment a limited transcribed speech corpus
- Multilingual acoustic model training
  - Using other-language speech to augment a limited transcribed speech corpus
- Removing reliance on pronunciation lexicons

   Grapheme based models and acoustically aided G2P
- Chain models
  - 10% more accurate STT, plus
  - 3x faster decoding, and 5x-10x faster training

# The Genesis of Chain Models

- Connectionist Temporal Classification
  - The latest shiny toy in neural network-based acoustic modeling for STT (ICASSP and InterSpeech 2015)
  - Nice STT improvements shown on Google datasets
  - We haven't seen STT gains on our datasets
- Chain Models
  - Inspired by (but quite different from) CTC
  - Sequence training of NNs without CE pre-training
  - Nice STT improvements over previous best systems
  - 3x decoding time speed-up; 5x-10x training speed-up

## 2006: A New Kid on the NNet Block

Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks

Alex Graves <sup>1</sup>	ALEX@IDSIA.CH		
Santiago Fernández <sup>1</sup>	SANTIAGO@IDSIA.CH		
Faustino Gomez <sup>1</sup>	TINO@IDSIA.CH		
Jürgen Schmidhuber <sup>1,2</sup>	JUERGEN@IDSIA.CH		
<sup>1</sup> Istituto Dalle Molle di Studi sull'Intelligenza Artificiale (II	OSIA), Galleria 2, 6928 Manno-Lugano, Switzerland		
<sup>2</sup> Technische Universität München (TUM), Boltzmannstr. 3, 85748 Garching, Munich, Germany			

#### Abstract

Many real-world sequence learning tasks require the prediction of sequences of labels from noisy, unsegmented input data. In speech recognition, for example, an acoustic signal is transcribed into words or sub-word units. Recurrent neural networks (RNNs) are powerful sequence learners that would seem well suited to such tasks. However, because they require pre-segmented training data, and post-processing to transform their outputs into label sequences, their applicability has so far been limited. This paper presents a novel method for training RNNs to label unsegmented sequences directly, thereby solving both problems. An experiment on the TIMIT speech corpus demonstrates its adbelling. While these approaches have proved successful for many problems, they have several drawbacks: (1) they usually require a significant amount of task specific knowledge, e.g. to design the state models for HMMs, or choose the input features for CRFs; (2) they require explicit (and often questionable) dependency assumptions to make inference tractable, e.g. the assumption that observations are independent for HMMs; (3) for standard HMMs, training is generative, even though sequence labelling is discriminative.

Recurrent neural networks (RNNs), on the other hand, require no prior knowledge of the data, beyond the choice of input and output representation. They can be trained discriminatively, and their internal state provides a powerful, general mechanism for modelling time series. In addition, they tend to be robust to temporal and spatial noise.

## 2015: The New Kid Comes of Age

### LEARNING ACOUSTIC FRAME LABELING FOR SPEECH RECOGNITION WITH RECURRENT NEURAL NETWORKS

Haşim Sak, Andrew Senior, Kanishka Rao, Ozan İrsoy, Alex Graves, Françoise Beaufays, Johan Schalkwyk

#### Google

{hasim, and rewsenior, kan ishkarao, gravesa, fsb, johans}@google.com

#### ABSTRACT

We explore alternative acoustic modeling techniques for large vocabulary speech recognition using Long Short-Term Memory recurrent neural networks. For an acoustic frame labeling task, we compare the conventional approach of cross-entropy (CE) training using fixed forced-alignments of frames and labels, with the Connectionist Temporal Classification (CTC) method proposed for labeling unsegmented sequence data. We demonstrate that the latter can be implemented with finite state transducers. We experiment with phones and context dependent HMM states as acoustic modeling units. We also investigate the effect of context in acoustic input by training unidirectional and bidirectional LSTM RNN models. We show that a bidirectional LSTM RNN CTC model using phone units can perform as well as an LSTM RNN model trained with CE using HMM state alignments. Finally, we also show the effect of sequence discriminative training on these models and show the first results for sMBR training of CTC models.

implemented in finite-state transducer (FST) framework and explain how CTC models can be used in decoding (Section 2.2). We also describe the use of sequence discriminative training with our sequence models (Section 2.3). In Section 4, we describe experiments with two acoustic modeling units – phones and HMM states. We also investigate the effect of acoustic context for LSTM RNN acoustic models by training unidirectional and bidirectional models.

#### 2. ACOUSTIC MODELING WITH LSTM RNN

There are a number of alternative approaches for acoustic modeling with neural networks for automatic speech recognition (ASR). Fundamentally the unit to be modeled by the network must be chosen (e.g. phone, HMM state, context dependent HMM state, diphone, word etc.). Training may use a hard (Viterbi) alignment with a single class label per frame, or a soft (Baum-Welch) alignment giving a probability distribution. Further, a variety of objective functions

## 2015: The New Kid Comes of Age

### ACOUSTIC MODELLING WITH CD-CTC-SMBR LSTM RNNS

Andrew Senior, Haşim Sak, Félix de Chaumont Quitry, Tara Sainath, Kanishka Rao

### Google

{hasim, and rewsenior, fcq, tsainath, kanishkarao}@google.com

#### ABSTRACT

This paper describes a series of experiments to extend the application of Context-Dependent (CD) long short-term memory (LSTM) recurrent neural networks (RNNs) trained with Connectionist Temporal Classification (CTC) and sMBR loss. Our experiments, on a noisy, reverberant voice search task, include training with alternative pronunciations and the application to child speech recognition; combination of multiple models, and convolutional input layers. We also investigate the latency of CTC models and show that constraining forward-backward alignment in training can reduce the delay for a real-time streaming speech recognition system. Finally we investigate transferring knowledge from one network to another through alignments.

Index Terms: Long Short Term Memory, Recurrent Neural Networks, Connectionist Temporal Classification, sequence discriminative training, knowledge transfer. the labels indicate the segmentation of the sequence with repeated labels indicating longer durations, with CTC an output may only be high for a single frame to indicate the presence of the symbol, with other frames labelled "blank," and duration information is discarded. During training CTC constantly aligns every sequence and trains to maximize the total probability of all valid label sequences. Because of the memory of the LSTM model this means that the outputs no longer need to occur at the same time as the input features to which they correspond.

In our previous work [2] we have shown that models with a blank symbol that are initialized with CTC can be improved upon with sMBR sequence-discriminative training. We then showed [3] that such models, using long-duration features (95ms of speech represented as 8 stacked overlapping log-mel filterbank features, generated with a 25ms window FFT every 10ms), downsampled and processed every 30ms, can outperform conventionally-trained LSTM models when using context dependent phone targets [5]. We use the term CD-CTC-sMBR LSTM RNN for these models.

## CTC, Explained ... in Pictures





## CTC, Explained ... in Pictures




### **DNN versus CTC: STT Performance**

Figures and Tables from Sak et al, ICASSP 2015



DNN	Target	CE	sMBR
LSTM	Senone	10.0%	8.9%
BLSTM	Senone	9.7%	9.1%





СТС	Target	CE	sMBR
LSTM	Phone	10.5%	9.4%
BLSTM	Phone	9.5%	8.5%
	CTC LSTM BLSTM	CTC Target LSTM Phone BLSTM Phone	CTCTargetCELSTMPhone10.5%BLSTMPhone9.5%

(i) CTC bidirectional LSTM with phone labels (10% PER)

### First, the Bad News ...

- We haven't been able to get CTC models to give us any noticeable improvement over our best (TDNN or LSTM-RNN) models on our data
  - It appears to be easier to get them to work when one has several 1000 hours of labeled speech
  - But we care about lower-resource scenarios

#### ... and then the Good News

- We are able to get similar improvements using a different model, which is inspired by ideas from the CTC papers
  - Use simple "1-state" HMMs for each CD phone
  - Reduce frame rate from 100 Hz to 33 Hz
  - Permit slack in the frame-to-state alignment

# Chain Models and LF-MMI Training

- A new class of acoustic models for hybrid STT
   "1-state" HMM for each context-dependent phone
  - LSTM/TDNNs compute state posterior probabilities
- MFCCs are down-sampled from 100Hz to 33Hz
   Inspired by CTC
- A new lattice-free MMI training method
  - Improved parallelization, sequence training on GPUs
    - Larger mini-batches, smaller I/O bandwidth
  - Does not require CE training before MMI training
  - Uses "flexible label alignment" inspired by CTC



# Lattice-Free MMI Training

- Denominator (phone) graph creation
  - Use a phone 4-gram language model, L
  - Compose H, C and L to obtain denominator graph
    - This FSA is the same for all utterances; suits GPU training
    - Use (heuristic) sentence-specific initial probabilities
- Numerator graph creation
  - Generate a phone graph using transcripts
    - This FSA encodes frame-by-frame alignment of HMM states
  - Permit some alignment "slack" for each frame/label
  - Intersect slackened FSA with the denominator FSA

# Lattice-free MMI Training (cont'd)

- LSTM-RNNs trained with this MMI training procedure are highly susceptible to over-fitting
- Essential to **regularize** the NN training process
  - A second output layer for CE training
  - Output L<sub>2</sub> regularization
  - Use a leaky HMM

Regularization		Hub-5 '00 Word Error Rate		
Cross Entropy	L <sub>2</sub> Norm	Leaky HMM	Total	SWBD
Ν	Ν	N	16.8%	11.1%
Y	Ν	Ν	15.9%	10.5%
N	Y	Ν	15.9%	10.4%
N	Ν	Y	16.4%	10.9%
Y	Y	Ν	15.7%	10.3%
Y	Ν	Y	15.7%	10.3%
N	Y	Y	15.8%	10.4%
Y	Y	Y	15.6%	10.4%

#### STT Results for Chain Models

300 hours of SWBD Training Speech; Hub-5 '00 Evaluation Set

Training Objective	Model (Size)	Total WER	SWBD WER
Cross-Entropy	TDNN A (16.6M)	18.2%	12.5%
CE + sMBR	TDNN A (16.6M)	16.9%	11.4%
	TDNN A (9.8M)	16.1%	10.7%
Lattice-free MMI	TDNN B (9.9M)	15.6%	10.4%
	TDNN C (11.2M)	15.5%	10.2%
LF-MMI + sMBR	TDNN C (11.2M)	15.1%	10.0%

- LF-MMI reduces WER by ca 10%-15% *relative*
- LF-MMI is better than standard CE + sMBR training (ca 8%)
- LF-MMI improves very slightly with additional sMBR training

#### Chain Models and LF-MMI Training

STT Performance on a Variety of Corpora

Corpus and Audio Type	Training Speech	CE + sMBR Error Rate	LF-MMI Error Rate
AMI IHM	80 hours	23.8%	22.4%
AMI SDM	80 hours	48.9%	46.1%
TED-LIUM	118 hours	11.3%	12.8%
Switchboard	300 hours	16.9%	15.5%
Fisher + SWBD	2100 hours	15.0%	13.3%

- Chain models with LF-MMI reduce WER by 6%-11% (relative)
- LF-MMI improves a bit further with additional sMBR training
- FL-MMI is 5x-10x faster to train, 3x faster to decode

# A Recap of Chain Models

- A new class of acoustic models for hybrid STT
  - "1-state" HMM for context-dependent phones
  - LSTM-RNN acoustic models (TDNN also compatible)
- A new lattice-free MMI training method
  - Better suited to using GPUs for parallelization
  - Does not require CE training before MMI training
- Improved speed and STT performance
  - 6%-8% relative WER reduction over previous best
  - 5-10x improvement in training time; 3x decoding time



Summary of Advanced Methods:

#### Staying Ahead in the STT Game

- STT technology is advancing very rapidly

   Amazon, Apple, Baidu, Facebook, Google, Microsoft
- Kaldi leads and keeps up with major innovations
  - From SGMMs to DNN (2012)
  - From "English" to low-resource languages (2013)
  - From CPUs to GPUs (2014)
  - From close-talking to far-field microphones (2015)
  - From well-curated to "wild type" corpora (2016)
  - Chain models for better STT, faster decoding (2017)
- and the list goes on ...

# Team Kaldi @ Johns Hopkins

- Sanjeev Khudanpur
- Daniel Povey
- Jan Trmal
- Guoguo Chen
- Pegah Ghahremani
- Vimal Manohar
- Vijayaditya Peddinti
- Hainan Xu
- Xiaohui Zhang
- ... and several others



# Kaldi Points-of-Contact

#### Kaldi mailing list

- <u>kaldi-help@googlegroups.com</u>
- Daniel Povey
  - <u>dpovey@gmail.com</u>
- Jan "Yenda" Trmal – <u>trmal@jhu.edu</u>
- Sanjeev Khudanpur
  - <u>khudanpur@jhu.edu</u>
  - 410-516-7024