

List of modifications made in HTS (for version 2.1)

Heiga ZEN

June 13, 2008

1 Modifications in Model Definition

In HTS, the HTK HMM definition (please see HTKBook [2] Chapter 7) has been modified to support MSD [3], stream-level tying, and adaptation of multi-stream HMMs. This section gives its brief description.

First, `<MSDInfo>` is added to global options of the HTK HMM definition language. The arguments to the `<MSDInfo>` option are the number of streams (default 1) and then for each stream, 0 (non-MSD stream) or 1 (MSD stream) of that stream. The full set of global options in HTS is given below.

```
globalOpts = option { option }
option      = <HmmSetId> string |
              <StreamInfo> short { short } |
              <MSDInfo> short { short } |
              <VecSize> short |
              <ProjSize> short |
              <InputXform> inputXform |
              <ParentXform> ~a macro |
              covkind |
              durkind |
              parmkind
```

Second, the number of mixture specification is modified to support stream-level tying structure as follows:

HTK	HTS
<code><State> 2</code>	<code><State> 2</code>
<code><NumMixes> 1 2</code>	
<code><SWeights> 2 0.9 1.1</code>	<code><SWeights> 2 0.9 1.1</code>
<code><Stream> 1</code>	<code><Stream> 1</code>
<code><Mixture> 1 1.0</code>	<code><NumMixes> 1</code>
<code><Mean> 4</code>	<code><Mixture> 1 1.0</code>
0.3 0.2 0.1 0.0	<code><Mean> 4</code>
<code><Variance> 4</code>	0.3 0.2 0.1 0.0
0.5 0.4 0.3 0.2	<code><Variance> 4</code>
<code><Stream> 2</code>	0.5 0.4 0.3 0.2
<code><Mixture> 1 0.4</code>	<code><Stream> 2</code>
<code><Mean> 2</code>	<code><NumMixes> 2</code>
1.0 2.0	<code><Mixture> 1 0.4</code>
<code><Variance> 2</code>	<code><Mean> 2</code>
4.0 8.0	1.0 2.0
<code><Mixture> 2 0.6</code>	<code><Variance> 2</code>
<code><Mean> 2</code>	4.0 8.0
2.0 9.0	<code><Mixture> 2 0.6</code>
<code><Variance> 2</code>	<code><Mean> 2</code>
3.0 6.0	2.0 9.0
	<code><Variance> 2</code>
	3.0 6.0

As you can see, `<NumMixes>` is moved from state-level to stream-level. This modification enables us to include the number of mixture component in the stream-level macro. Based on this implementation, stream-level macro was added. The various distinct points in the hierarchy of HMM parameters which can be tied in HTS is as follows:

```

~s  shared state distribution
~p  shared stream
~m  shared Gaussian mixture component
~u  shared mean vector
~v  shared diagonal variance vector
~i  shared inverse full covariance matrix
~c  shared Cholesky U matrix
~x  shared arbitrary transform matrix
~t  shared transition matrix
~d  shared duration parameters
~w  shared stream weight vector

```

Note that the `~p` macro is used by the HMM editor HHed for building tied mixture systems in the original HTK macro definition.

The resultant state definition of in the modified HTK HMM definition language is as follows:

```

state      = <State> short stateinfo
stateinfo  = ~s macro |
             [ weights ] stream { stream } [ duration ]
macro      = string
weights    = ~w macro | <SWeights> short vector
vector     = float { float }
stream     = [ <Stream> short ] streaminfo
streaminfo = ~p macro | [ <Stream> short ] [mixes] (mixture { mixture } | tmixpdf | discpdf)
mixes      = <NumMixes> short {short}
tmixpdf    = <TMix> macro weightList
weightList = repShort { repShort }
repShort   = short [ * char ]
discpdf    = <DProb> weightList
mixture    = [ <Mixture> short float ] mixpdf
mixpdf     = ~m macro | mean cov [ <GConst> float ]
mean       = ~u macro | <Mean> short vector
cov        = var | inv | xform
var        = ~v macro | <Variance> short vector
inv        = ~i macro |
             (<InvCovar> | <LLTCovar>) short tmatrix
xform      = ~x macro | <Xform> short short matrix
matrix     = float {float}
tmatrix    = matrix

```

It should be noted that `<Stream>` can doubly be specified in both stream and streaminfo. This is because `<Stream>` in `~p` macro is essential to specify stream index of this macro. This stream index information is used in various HTS functions to check stream consistency.

Third, to support multi-stream HMM adaptation, the HTK HMM definition language for baseclasses is modified. A baseclass is defined as

```

baseClass  = ~b macro baseopts classes
baseopts   = <MMFIdMask> string <Parameters> baseKind [<StreamInfo>] <NumClasses> int
StreamInfo = short { short } |
baseKind   = MIXBASE | MEANBASE | COVBASE
classes    = <Class> int itemlist { classes }

```

where `<StreamInfo>` is optionally added to specify the stream structure.

2 Added Configuration Variables

A number of configuration variables have been added to HTK to control new functions implemented in HTS. Their names, default values, and brief descriptions are as follows:

Module	Name	Default	Description
HADAPT	SAVEFULLC	F	Save transformed model set in full covariance form
	USESMAP	F	Use structural MAP criterion [4]
	SMAPSIGMA	1.0	Prior parameter for SMAP criterion
	SAVEALLSMAPXFORM	T	Save all (unnecessary) linear transforms estimated in SMAPLR/CSMAPLR
	BANDWIDTH		Bandwidth of transformation matrices [5]
	DURUSEBIAS	F	Specify a bias with linear transforms
	DURSPLITTHRESH	1000.0	Minimum occupancy to generate a transform for state duration model set
	DURTRANSKIND	MLLRMEAN	Transformation kind
	DURBLOCKSIZE	full	Block structure of transform for state duration model set
	DURBANDWIDTH		Bandwidth of transformation matrices for state duration model set
	DURBASECLASS	global	Macroname of baseclass for state duration model set
	DURREGTREE		Macroname of regression tree for state duration model set
	DURADAPTKIND	BASE	Use regression tree or base classes to adapt state duration model set
HFB	MAXSTDDEVCOEF	10	Maximum duration to be evaluated
	MINDUR	5	Minimum duration to be evaluated
HMAP	APPLYVFLOOR	T	Apply variance floor to model set
HGEN	MAXEMITER	20	Maximum # of EM iterations
	EMEPSILON	1.0E-4	Convergence factor for EM iteration
	RNDPARMEAN	0.0	Mean of Gaussian noise for random generation [6]
	RNDPARVAR	1.0	Variance of Gaussian noise for random generation
	USEGV	F	Use speech parameter generation algorithm considering GV [7]

Module	Name	Default	Description
	CDGV	F	Use context-dependent GV model set
	LOGGV	F	Use logarithmic GV instead of linear GV
	MAXGVITER	F	Max iterations in the speech parameter generation considering GV
	GVEPSILON	1.0E-4	Convergence factor for GV iteration
	MINEUCNORM	1.0E-2	Minimum Euclid norm of a gradient vector
	STEPINIT	1.0	Initial step size
	STEPDEC	0.5	Step size deceleration factor
	STEPINC	1.2	Step size acceleration factor
	HMMWEIGHT	1.0	Weight for HMM output prob
	GVWEIGHT	1.0	Weight for GV output prob
	OPTKIND	NEWTON	Optimization method
	RNDFLAGS		Random generation flag
	GVMODELMMF		GV MMF file
	GVHMMLIST		GV model list
	GVMODELDIR		Dir containing GV models
	GVMODELEXT		Ext to be used with above Dir
	GVOFFMODEL		Model names to be excluded from GV calculation
HMODEL	IGNOREVALUE	-1.0E+10	Ignore value to indicate zero-dimensional space in multi-space probability distribution
HCOMPV	NSHOWELEM	12	# of vector elements to be shows
	VFLOORSSCALE	0.0	variance flooring scale
	VFLOORSSCALESTR		variance flooring scale vector for streams
HEREST	APPLYVFLOOR	T	Apply variance floor to model set
	DURMINVAR	0.0	Minimum variance floor for state duration model set
	DURVARFLOORPERCENTILE	0	Maximum number of Gaussian components (as the percentage of the total Gaussian components in the system) to undergo variance floor for state duration model set
	APPLYDURVARFLOOR	T	Apply variance floor to state duration model set
	DURMAPTAU	0.0	MAP tau for state duration model set [8]
	ALIGNDURMMF		State duration MMF file for alignment (2-model reest)

Module	Name	Default	Description
	ALIGNDURLIST		State duration model list for alignment (2-model reest)
	ALIGNDURDIR		Dir containing state duration models for alignment (2-model reest)
	ALIGNDUREXT		Ext to be used with above Dir (2-model reest)
	ALIGNDURXFORMEXT		Input transform ext for state duration model set to be used with 2-model reest
	ALIGNDURXFORMDIR		Input transform dir for state duration model set to be used with 2-model reest
	DURINXFORMMASK		Input transform mask for state duration model set (default output transform mask)
	DURPAXFORMMASK		Parent transform mask for state duration model set (default output parent mask)
HHed	USEPATTERN	F	Use pattern instead of base phone for tree-based clustering
	SINGLETREE	F	Construct single tree for each state position
	APPLYMDL	F	Use the MDL criterion for tree-based clustering [9]
	IGNORESTRW	F	Ignore stream weight in tree-based clustering
	REDUCEMEM	F	Use reduced memory implementation of tree-based clustering
	MINVAR	1.0E-6	Minimum variance floor for model set
	MDLFACTOR	1.0	Factor to control the model complexity term in the MDL criterion
	MINLEAFOCC	0.0	Minimum occupancy count in each leaf node
	MINMIXOCC	0.0	Minimum occupancy count in each mixture component
	SHRINKOCCTHRESH		Minimum occupancy count in decision trees shrinking
HMGENS	SAVEBINARY	F	Save generated parameters in binary
	OUTPDF	F	Output pdf sequences
	PARMGENTYPE	0	Type of parameter generation algorithm [10]
	MODELALIGN	F	Use model-level alignments given from label files to determine model-level durations

Module	Name	Default	Description
	STATEALIGN	F	Use state-level alignments given from label files to determine state-level durations
	USEALIGN	F	Use model-level alignments to prune EM-based parameter generation algorithm
	USEHMMFB	F	Do not use state duration models in the EM-based parameter generation algorithm
	INXFORMMASK		Input transform mask
	PAXFORMMASK		Parent transform mask
	PDFSTRSIZE		# of PdfStreams
	PDFSTRORDER		Size of static feature in each PdfStream
	PDFSTREXT		Ext to be used for generated parameters from each PdfStream
	WINEXT		Ext to be used for window coefficients file
	WINDIR		Dir containing window coefficient files
	WINFN		Name of window coefficient files

Other configuration variables in HTK can also be used with HTS. Please refer to HTKBook [2] Chapter 18 for others.

3 Added Command-Line Options

Various new command-line options have also been added to HTK tools. They are listed as follows:

HInit

Option		Default
-g	Ignore outlier vector in MSD	on

HRest

Option		Default
-g s	output duration model to file s	none
-o fn	Store new hmm def in fn (name only)	outDir/srcfn

HERest

Option		Default
-b	use an input linear transform for dur models	off
-f s	extension for new duration model files	as src
-g s	output duration model to file s	none
-n s	dir to find duration model definitions	current
-q s	save all xforms for duration to TMF file s	TMF
-u tmvwapd	update t)rans m)eans v)ars w)ghts a)daptation xform p)rrior used s)semi-tied xform d) switch to duration model update flag	tmvw
-y s	extension for duration model files	none
-N mmf	load duration macro file mmf	
-R dir	dir to write duration macro files	current
-W s [s]	set dir for duration parent xform to s and optional extension	off
-Y s [s]	set dir for duration input xform to s and optional extension	none
-Z s [s]	set dir for duration output xform to s	none

HHed

Option		Default
-a f	factor to control the second term in the MDL	1.0
-i	ignore stream weight	off
-m	apply MDL principle for clustering	off
-p	use pattern instead of base phone	off
-r	reduce memory usage on clustering	off
-s	construct single tree	off
-v f	Set minimum variance to f	1.0E-6

HMGenS

Option		Default
--------	--	---------

-a	Use an input linear transform for HMMs	off
-b	Use an input linear transform for dur models	off
-c n	type of parameter generation algorithm	0
	0: both mix and state sequences are given	
	1: state sequence is given, but mix sequence is hidden	
	2: both state and mix sequences are hidden	
-d s	dir to find hmm definitions	current
-e	use model alignment from label for pruning	off
-f f	frame shift in 100 ns	50000
-g f	Mixture pruning threshold	10.0
-h s [s]	set speaker name pattern to s, optionally set parent patterns	*.%%%
-m	use model alignment for duration	off
-n s	dir to find duration model definitions	current
-p	output pdf sequences	off
-r f	speaking rate factor (f<1: fast f>1: slow)	1.0
-s	use state alignment for duration	off
-t f [i l]	set pruning to f [inc limit]	inf
-v f	threshold for switching spaces for MSD	0.5
-x s	extension for hmm files	none
-y s	extension for duration model files	none
-E s [s]	set dir for parent xform to s and optional extension	off
-G fmt	Set source label format to fmt	as config
-H mmf	Load HMM macro file mmf	
-I mlf	Load master label file mlf	
-J s [s]	set dir for input xform to s and optional extension	none
-L dir	Set input label (or net) dir	current
-M dir	Dir to write HMM macro files	current
-N mmf	Load duration macro file mmf	
-S f	Set script file to f	none
-T N	Set trace flags to N	0
-V	Print version information	off
-W s [s]	set dir for duration parent xform to s and optional extension	off
-X ext	Set input label (or net) file ext	lab
-Y s [s]	set dir for duration input xform to s and optional extension	none

Please also refer to HTKBook [2] Chapter 17 for other command-line options.

4 Added Commands and Modifications in HHed

Some HHed commands have been added in HTS. They are as follows:

AX filename	- Set the Adapt XForm to filename
CM directory	- Convert models to pdf for speech synthesizer
CT directory	- Convert trees/questions for speech synthesizer
DM type macroname	- Delete macro from model-set
DR id	- Convert decision trees to a regression tree
DV	- Convert full covariance to diagonal variances
IT filename	- Clustering while imposing loaded tree structure If any empty leaf nodes exist, loaded trees are pruned and then saved to filename
IX filename	- Set the Input Xform to filename
PX filename	- Set the Parent Xform to filename
// comment	- Comment line (ignored)

In many HHed commands, we are required to specify item lists to specify a set of items to be processed. In HTS, item list specification has been modified to specify stream-level items.

```
itemList  = "{ " itemSet { " , " itemSet } " }"
itemSet   = hmmName . [ "transP" | "state" state ]
hmmName=  ident | identList
identList = "( " ident { " , " ident } " )"
ident     = < char | metachar >
metachar  = "?" | "*"
state     = index [ " ." stateComp ]
index     = "[ " intRange { " , " intRange } "]"
intRange  = integer [ " ." integer ]
stateComp = "dur" | "weights" | stream
stream    = [ " stream" index ] [ " .mix" mix ]
mix       = index [ " ." ( "mean" | "cov" ) ]
```

For example,

```
TI str1 { *.state[2].stream[1]}
```

denotes tying streams in state 2 of all phonemes.

Appendix A History of HTS

- **Version 1.0 (December 2002)**
 - Based on HTK-3.2.
 - Tree-based clustering based on the MDL criterion [9].
 - Stream-dependent tree-based clustering [11].
 - Multi-space probability distributions (MSD) [3].
 - State duration modeling and clustering [12].
 - Speech parameter generation algorithm [10].
 - Demo using the CMU Communicator database.
- **Version 1.1 (May 2003)**
 - Based on HTK-3.2.
 - Small run-time synthesis engine.
 - Demo using the CSTR TIMIT database.
 - HTS voices for the Festival speech synthesis system [13].
- **Version 1.1.1 (December 2003)**
 - Based on HTK-3.2.1.
 - Variance flooring for MSD-HMMs.
 - Post-filtering [14].
 - Demo using the CMU ARCTIC database.
 - Demo using the Nitech Japanese database.
 - HTS voice for the Galatea toolkit [15].
- **Version 2.0 (December 2006) [16]**
 - Based on HTK-3.4.
 - Support generating state duration PDFs in HRest.
 - Phoneme boundaries can be given to HERest using the -e option.
 - Reduced-memory implementation of tree-based clustering in HHed with the -r option.
 - Each decision tree can have a name with regular expressions in HHed with the -p option.
 - Flexible model structures in HMGenS.
 - Speech parameter generation algorithm based on the EM algorithm [10] in HMGenS.
 - Random generation algorithm [6] in HMGenS.
 - State or phoneme-level alignments can be given to HMGenS.
 - The interface of HMGenS has been switched to HERest-style.
 - Various kinds of linear transformations for MSD-HMMs are supported in HERest.
 - * Constrained MLLR based adaptation [17].
 - * Adaptive training based on constrained MLLR [17].
 - * Precision matrix modeling based on semi-tied covariance matrices [18].
 - * Heteroscedastic linear discriminant analysis (HLDA) based feature transform [19].
 - * Phonetic decision trees can be used to define regression classes for adaptation [20]
 - * Adapted HMMs can be converted to the run-time synthesis engine format.
 - Maximum a posteriori (MAP) adaptation [8] for MSD-HMMs in HERest.
- **Version 2.0.1 (May 2007)**
 - Based on HTK-3.4.
 - Band structure for linear transforms [5].
 - Speaker interpolation [21].
 - Stream-dependent variance flooring scales.
 - Demo scripts support LSP-type spectral parameters.
 - β version of the runtime synthesis engine API.
- **Version 2.1 (June 2008)**
 - Based on HTK-3.4.
 - Released under the New and Simplified BSD license (<http://www.opensource.org/>).

- Include simple documentation
- Support 64-bit compile
- HERest supports HSMM training and adaptation [22, 23].
- HAdapt supports CSMAPLR adaptation [4].
- HMGenS supports speech parameter generation from HSMMs.
- Speech parameter generation algorithm considering GV [7].
- Random generation of transitions, durations, and mixture components in HMGenS.
- Add DM command to HHed to delete an existing macro from MMF.
- Add IT command to HHed to impose pre-constructed trees in clustering.
- HHed MU command supports '*2' style mixing up.
- HHed MU command supports mixture-level occupancy threshold in mixing up.
- First stable version of the runtime synthesis engine API (hts_engine API).

References

- [1] K. Tokuda, H. Zen, J. Yamagishi, A.W. Black, T. Masuko, S. Sako, T. Toda, T. Nose, and K. Oura. The HMM-based speech synthesis system (HTS). <http://hts.sp.nitech.ac.jp/>.
- [2] S. Young, G. Evermann, M. Gales, T. Hain, D. Kershaw, X.-Y. Liu, G. Moore, J. Odell, D. Ollason, D. Povey, V. Valtchev, and P. Woodland. *The Hidden Markov Model Toolkit (HTK) version 3.4*, 2006. <http://htk.eng.cam.ac.uk/>.
- [3] K. Tokuda, T. Masuko, N. Miyazaki, and T. Kobayashi. Multi-space probability distribution HMM. *IEICE Trans. Inf. & Syst.*, E85-D(3):455–464, Mar. 2002.
- [4] Y. Nakano, M. Tachibana, J. Yamagishi, and T. Kobayashi. Constrained structural maximum a posteriori linear regression for average-voice-based speech synthesis. In *Proc. Interspeech*, pages 2286–2289, 2006.
- [5] L. Qin, Y.-J. Wu, Z.-H. Ling, and R.-H. Wang. Improving the performance of HMM-based voice conversion using context clustering decision tree and appropriate regression matrix. In *Proc. of Interspeech (ICSLP)*, pages 2250–2253, 2006.
- [6] K. Tokuda, H. Zen, and T. Kitamura. Reformulating the HMM as a trajectory model. In *Proc. Beyond HMM – Workshop on statistical modeling approach for speech recognition*, 2004.
- [7] T. Toda and K. Tokuda. A speech parameter generation algorithm considering global variance for HMM-based speech synthesis. *IEICE Trans. Inf. & Syst.*, E90-D(5):816–824, 2007.
- [8] J.L. Gauvain and C.-H. Lee. Maximum a posteriori estimation for multivariate Gaussian mixture observations of Markov chains. *IEEE Trans. on Speech & Audio Process.*, 2(2):291–298, 1994.
- [9] K. Shinoda and T. Watanabe. MDL-based context-dependent subword modeling for speech recognition. *J. Acoust. Soc. Jpn.(E)*, 21(2):79–86, 2000.
- [10] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura. Speech parameter generation algorithms for HMM-based speech synthesis. In *Proc. ICASSP*, pages 1315–1318, 2000.
- [11] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura. Simultaneous modeling of spectrum, pitch and duration in HMM-based speech synthesis. In *Proc. Eurospeech*, pages 2347–2350, 1999.
- [12] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura. Duration modeling for HMM-based speech synthesis. In *Proc. ICSLP*, pages 29–32, 1998.
- [13] A.W. Black, P. Taylor, and R. Caley. The festival speech synthesis system. <http://www.festvox.org/festival/>.
- [14] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura. Incorporation of mixed excitation model and postfilter into HMM-based text-to-speech synthesis. *IEICE Trans. Inf. & Syst. (Japanese Edition)*, J87-D-II(8):1563–1571, Aug. 2004.
- [15] Galatea – An open-source toolkit for anthropomorphic spoken dialogue agent. <http://hil.t.u-tokyo.ac.jp/galatea/>.
- [16] H. Zen, T. Nose, J. Yamagishi, S. Sako, T. Masuko, A.W. Black, and K. Tokuda. The HMM-based speech synthesis system version 2.0. In *Proc. ISCA SSW6*, pages 294–299, 2007.
- [17] M.J.F. Gales. Maximum likelihood linear transformations for HMM-based speech recognition. *Computer Speech & Language*, 12(2):75–98, 1998.
- [18] M.J.F. Gales. Semi-tied covariance matrices for hidden Markov models. *IEEE Transactions on Speech and Audio Processing*, 7(3):272–281, 1999.
- [19] M.J.F. Gales. Maximum likelihood multiple projection schemes for hidden Markov models. *IEEE Trans. Speech & Audio Process.*, 10(2):37–47, 2002.
- [20] J. Yamagishi, M. Tachibana, T. Masuko, and T. Kobayashi. Speaking style adaptation using context clustering decision tree for HMM-based speech synthesis. In *Proc. ICASSP*, pages 5–8, 2004.
- [21] T. Yoshimura, T. Masuko, K. Tokuda, T. Kobayashi, and T. Kitamura. Speaker interpolation for HMM-based speech synthesis system. *J. Acoust. Soc. Jpn. (E)*, 21(4):199–206, 2000.
- [22] H. Zen, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura. A hidden semi-Markov model-based speech synthesis system. *IEICE Trans. Inf. & Syst.*, E90-D(5):825–834, 2007.
- [23] J. Yamagishi and T. Kobayashi. Average-voice-based speech synthesis using HSMM-based speaker adaptation and adaptive training. *IEICE Trans. Inf. & Syst.*, E90-D(2):533–543, 2007.