

# List of modifications made in HTS (for version 2.3 alpha)

HTS Working Group

December 25, 2012

## 1 Modifications in Model Definition

In HTS, the HTK HMM definition (please see HTKBook [1] Chapter 7) has been modified to support MSD [2], 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>&lt;State&gt; 2</code>	<code>&lt;State&gt; 2</code>
<code>&lt;NumMixes&gt; 1 2</code>	
<code>&lt;SWeights&gt; 2 0.9 1.1</code>	<code>&lt;SWeights&gt; 2 0.9 1.1</code>
<code>&lt;Stream&gt; 1</code>	<code>&lt;Stream&gt; 1</code>
	<code>&lt;NumMixes&gt; 1</code>
<code>&lt;Mixture&gt; 1 1.0</code>	<code>&lt;Mixture&gt; 1 1.0</code>
<code>&lt;Mean&gt; 4</code>	<code>&lt;Mean&gt; 4</code>
<code>0.3 0.2 0.1 0.0</code>	<code>0.3 0.2 0.1 0.0</code>
<code>&lt;Variance&gt; 4</code>	<code>&lt;Variance&gt; 4</code>
<code>0.5 0.4 0.3 0.2</code>	<code>0.5 0.4 0.3 0.2</code>
<code>&lt;Stream&gt; 2</code>	<code>&lt;Stream&gt; 2</code>
	<code>&lt;NumMixes&gt; 2</code>
<code>&lt;Mixture&gt; 1 0.4</code>	<code>&lt;Mixture&gt; 1 0.4</code>
<code>&lt;Mean&gt; 2</code>	<code>&lt;Mean&gt; 2</code>
<code>1.0 2.0</code>	<code>1.0 2.0</code>
<code>&lt;Variance&gt; 2</code>	<code>&lt;Variance&gt; 2</code>
<code>4.0 8.0</code>	<code>4.0 8.0</code>
<code>&lt;Mixture&gt; 2 0.6</code>	<code>&lt;Mixture&gt; 2 0.6</code>
<code>&lt;Mean&gt; 2</code>	<code>&lt;Mean&gt; 2</code>
<code>2.0 9.0</code>	<code>2.0 9.0</code>
<code>&lt;Variance&gt; 2</code>	<code>&lt;Variance&gt; 2</code>
<code>3.0 6.0</code>	<code>3.0 6.0</code>

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
	USEMAPLR	F	Use MAP criterion [3]
	USEVBLR	F	Use VB criterion [4]
	USESTRUCTURALPRIOR	F	Perform structural approach (e.g., SMAPLR, CSMAPLR, and SVBLR)
	PRIORSCALE	1.0	Prior parameter for MAPLR and VBLR
	SAVEALLNODEXFORMS	T	Save all (unnecessary) linear transforms estimated in SMAPLR/CSMAPLR/SVBLR
	BANDWIDTH		Bandwidth of transformation matrices [5]
	DURUSEBIAS	F	Specify a bias with linear transforms
	DURSPPLITTHRESH	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 number of EM iterations
	EMEPSILON	1.0E-4	Convergence factor for EM iteration

Module	Name	Default	Description
	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]
	USEGVPST		GV calculation flag for each stream
	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
	GVINITWEIGHT	1.0	Initial weight of GV
	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
	USEDAEM	F	Use the DAEM-based parameter generation algorithm
	DAEMITER	20	Number of iterations for the DAEM-based parameter generation algorithm
	DAEMTEMPSCHEDULE	1.0	Temperature schedule parameter for DAEM
HMODEL	IGNOREVALUE	-1.0E+10	Ignore value to indicate zero-dimensional space in multi-space probability distribution
HCOMPV	NSHOWELEM	12	Number of vector elements to be shows
	VFLOORSCALE	0.0	variance flooring scale
	VFLOORSCALESTR		variance flooring scale vector for streams

Module	Name	Default	Description
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)
	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

Module	Name	Default	Description
	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
	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		Number 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
HSMMALIGN	INXFORMMASK		Input transform mask
	PAXFORMMASK		Parent transform mask
	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)

Module	Name	Default	Description
HMGETOOL	SAVEBINARY	F	Save estimated models in binary
	PDFSTRSIZE		Number of PdfStreams
	PDFSTRORDER		Size of static feature in each Pdf-Stream
	MGETRNFFLAG		Whether perform MGE train for the stream
	GVTRNFFLAG		Whether incorporate GV component into MGE training
	ACCERRFLAG		Accumulate generation error
	INVQUASIZE		Bandwidth of quasi-diagonal inversion matrix
	VARWINSIZE		Window for local variance calculation
	GVDISTWGHT		GV weights
	GAINWTFLAG		gain weight for generation error
	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 [1] 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
-k f	set temperature parameter for DAEM training	1.0
-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
-q n	use reference tree for clustering	off
	0: clustering is stopped by threshold	
	1: clustering is stopped when leaf don't have occ	
	2: clustering is stopped by threshold	
	After that, standard clustering is not performed	
	3: clustering is stopped when leaf don't have occ	



After that, standard clustering is not performed

-r n	reduce memory usage on clustering	0
	0: no memory reduction	
	1: mid reduction but fast	
	2: large reduction but slow	
-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
-A	Print command line arguments	off
-B	Save HMMs/transforms as binary	off
-C cf	Set config file to cf	default
-D	Display configuration variables	off
-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

## HSMMALIGN

Option		Default
-a	Use an input linear transform for HMMs	off
-b	Use an input linear transform for dur models	off
-c	Prune by time information of label	off
-d s	Dir to find hmm definitions	current
-f	Output full state alignment	off
-h s [s]	Set speaker name pattern to s, optionally set parent patterns	*.%%%
-n s	Dir to find duration model definitions	current
-m dir	Set output label dir	current
-r ext	Output label file extension	lab
-s s	print statistics to file s	off
-t i	Set pruning threshold	off
-w f	Duration weight	1.0
-x s	Extension for hmm files	none
-y s	Extension for duration model files	none
-A	Print command line arguments	off
-C cf	Set config file to cf	default
-D	Display configuration variables	off
-E s [s]	set dir for parent xform to s and optional extension	off
-F fmt	Set source data format to fmt	as config
-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
-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

## HMGETOOL

Option		Default
-a i j	i: max times to shift the boundary j: max length for each boundary shifting	none 1
-b i j	i: end iteration for boundary adjustment j: window size for boundary adjustment	none 5
-c	output the process data	off

-d dir	HMM definition directory	none
-e	enable limit the updating rate for each step	off
-f r	frame rate	50000
-g	enable multiply variance ratio for mean updating	off
-i i j	start/end iteration index of MGE training	0 0
-j flg	0: eval 1: train 2: adapt	1
-l dir	output label directory	none
-o ext	HMM def file extension	none
-p a b	parameter for step size: $1/(a + b*n)$	1000.0 1.0
-r file	load HMM for reference	none
-s file	updating scale file	none
-u mvwa	update t)rans m)eans v)ars w)ghts a)daptation xform	none
-v f	threshold for switching spaces for MSD	0.5
-w f	distance weight for gv component	1.0
-x ext	label file extension	lab
-A	Print command line arguments	off
-B	Save HMMs/transforms as binary	off
-C cf	Set config file to cf	default
-D	Display configuration variables	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
-K s [s]	set dir for output 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
-S f	Set script file to f	none
-T N	Set trace flags to N	0
-V	Print version information	off
-X ext	Set input label (or net) file ext	lab

Please also refer to HTKBook [\[1\]](#) 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
AX filename state_mapping_table mmf list output_filename
                        - Set the Adapt XForm to another model and save one
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
JM hmmFile itemlist   - Join Models on stream or state level
PX filename           - Set the Parent Xform to filename
SM smtable mmf hlist  - Output KLD-based state mapping table
// 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.
  - HHED supports tree-based clustering based on the MDL criterion [9].
  - HHED supports stream-dependent tree-based clustering [11].
  - HMODEL supports multi-space probability distributions (MSD) [2].
  - HEREST can generate state duration modeling [12].
  - Speech parameter generation algorithm [10] is implemented in HGEN and HMGENS.
  - Demo using the CMU Communicator database.
- **Version 1.1 (May 2003)**
  - Based on HTK-3.2.
  - Small run-time synthesis engine (hts\_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.
  - HCOMPV supports variance flooring for MSD-HMMs.
  - Demo using the CMU ARCTIC database [14].
  - Demo using the Nitech Japanese database.
  - Demo supports post-filtering [15].
  - HTS voice for the Galatea toolkit [16].
- **Version 2.0 (December 2006) [17]**
  - Based on HTK-3.4.
  - Support generating state duration PDFs in HREST.
  - Phoneme boundaries can be given to HEREST using the -e option [18].
  - 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 [6].
  - 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 HADAPT.
    - \* Constrained MLLR based adaptation [19].
    - \* Adaptive training based on constrained MLLR [19].
    - \* Precision matrix modeling based on semi-tied covariance matrices [20].
    - \* Heteroscedastic linear discriminant analysis (HLDA) based feature transform [21].
    - \* Phonetic decision trees can be used to define regression classes for adaptation [22].
    - \* Adapted HMMs can be converted to the run-time synthesis engine format.
  - Maximum a posteriori (MAP) adaptation [8] for MSD-HMMs in HMAP.
- **Version 2.0.1 (May 2007)**
  - Based on HTK-3.4.
  - HADAPT supports band structure for linear transforms [5].
  - HCOMPV supports stream-dependent variance flooring scales.
  - Demo support LSP-type spectral parameters.
  - $\beta$  version of the runtime synthesis engine API (hts\_engine API).
  - hts\_engine API supports speaker interpolation [23].
- **Version 2.1 (June 2008)**
  - Based on HTK-3.4.
  - Released under the Modified BSD license [24].

- Simple documentation.
- 64-bit compile support.
- MAXSTRLEN (maximum length of strings), MAXFNAMELEN (maximum length of filenames), PAT.LEN (maximum length of pattern strings), and SMAX (maximum number of streams) defined in HShell.h can be set through configure script.
- HFB supports the forward/backward algorithm for hidden semi-Markov models (HSMMs) [25, 26].
- HADAPT supports SMAPLR/CSMAPLR adaptation [27, 3].
- HGEN supports speech parameter generation algorithm considering global variance (GV) [7].
- HGEN supports random generation of transitions, durations, and mixture components.
- HEREST supports HSMM training and adaptation.
- HMGENS supports speech parameter generation from HSMMs.
- Add DM command to HHED to delete an existing macro from MMF.
- Add IT command to HHED to impose pre-constructed trees in clustering.
- Add JM command to HHED to join models on state or stream level.
- 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).
- **Version 2.1.1 (May 2010)**
  - Based on HTK-3.4.1
  - WFST converter for forced-alignment of HSMM is implemented in HFST [28].
  - Demo support context-dependent GV without silent and pause phoneme.
  - Add initial GV weight for parameter generation in HMGENS.
  - Add memory reduction options for context-clustering in HHED.
  - Add model-level alignments given from label of singing voice to determine note-level durations.
  - Demo using the Nitech Japanese database for singing voice synthesis [29].
  - The API of runtime synthesis engine, hts\_engine API, is splitted from HTS itself and moved to SourceForge.
- **Version 2.2 (July 2011)**
  - Support DAEM algorithm in parameter estimation step [30].
  - Support KLD-based state-mapping and cross-lingual speaker adaptation in HHED [31].
  - Add stand-alone HSMM-based forced-alignment command, HSMMALIGN, instead of HFST [28].
  - Add HMGETOOL for MGE training [32].
  - Context-clustering can be started in the middle of the tree building.
  - Change sampling frequency of demo from 16kHz to 48kHz.
  - Support bark critical-band based aperiodic measure.
  - Change speakers of Brazilian Portuguese and Japanese song demo.
  - Release slides as a tutorial of HMM-based speech synthesis.
- **Version 2.3 alpha (December 2012)**
  - Add VBLR adaptation in HEREST [4].
  - Add DAEM-based parameter generation algorithm in HMGENS.
  - Support DP search to determine state duration when the model alignments are given in HMGENS.
  - Speed up context-clustering by calculating differences between answers to current and previous questions.
  - Add LSP postfilter in demo [33].
  - Support mel-cepstrum based aperiodic measure in STRAIGHT demo.
  - Support new HTS voice format for hts\_engine API.
  - Turn off spectrum normalization in STRAIGHT demo.

## References

- [1] 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/>.
- [2] 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.
- [3] 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.
- [4] K. Yu and M. Gales. Bayesian adaptation inference and adaptive training. *IEEE Trans. Audio, Speech and Language Process.*, 15:1932–1943, 2007.
- [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] J. Kominek and A.W. Black. CMU ARCTIC databases for speech synthesis. Technical Report CMU-LTI-03-177, Carnegie Mellon University, 2003.
- [15] 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.
- [16] Galatea – An open-source toolkit for anthropomorphic spoken dialogue agent. <http://hil.t.u-tokyo.ac.jp/galatea/>.
- [17] 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.
- [18] D. Huggins-Daines and A. Rudnicky. A constrained Baum-Welch algorithm for improved phoneme segmentation and efficient training. In *Proc. of Interspeech*, pages 1205–1208, 2006.
- [19] M.J.F. Gales. Maximum likelihood linear transformations for HMM-based speech recognition. *Computer Speech & Language*, 12(2):75–98, 1998.
- [20] M.J.F. Gales. Semi-tied covariance matrices for hidden Markov models. *IEEE Transactions on Speech and Audio Processing*, 7(3):272–281, 1999.
- [21] M.J.F. Gales. Maximum likelihood multiple projection schemes for hidden Markov models. *IEEE Trans. Speech & Audio Process.*, 10(2):37–47, 2002.
- [22] 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.
- [23] 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.
- [24] <http://www.opensource.org/licenses/category>.
- [25] 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.
- [26] 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.
- [27] O. Shiohan, Y. Myrvoll, and C.-H. Lee. Structural maximum a posteriori linear regression for fast HMM adaptation. *Computer Speech & Language*, 16(3):5–24, 2002.
- [28] K. Oura, H. Zen, Y. Nankaku, A. Lee, and K. Tokuda. A fully consistent hidden semi-Markov model-based speech recognition system. *IEICE Trans. Inf. & Syst.*, E91-D(11):2693–2700, 2008.
- [29] K. Oura, A. Mase, T. Yamada, S. Muto, Y. Nankaku, and K. Tokuda. Recent development of the HMM-based singing voice synthesis system — Sinsy. In *Proc. ISCA SSW7*, pages 211–216, 2010.
- [30] Y. Itaya, H. Zen, Y. Nankaku, C. Miyajima, K. Tokuda, and T. Kitamura. Deterministic annealing EM algorithm in acoustic modeling for speaker and speech recognition. *IEICE Trans. Inf. & Syst.*, E88-D(3):425–431, 2005.
- [31] Y.-J. Wu and K. Tokuda. State mapping based method method for cross-lingual speaker adaptation in HMM-based speech synthesis. In *Proc. of Interspeech (ICSLP)*, pages 420–423, 2009.
- [32] Y.-J. Wu and R.-H. Wang. Minimum generation error training for HMM-based speech synthesis. In *Proc. ICASSP*, pages 89–92, 2006.
- [33] Z.-H. Ling, Y.-J. Wu, Y.-P. Wang, L. Qin, and R.-H. Wang. USTC system for Blizzard Challenge 2006 an improved HMM-based speech synthesis method. In *Blizzard Challenge Workshop*, 2006.