

# VidTIMIT Dataset Documentation

Conrad Sanderson  
conradsand at ieee dot org

Version: 4 June 2008

## License

The VidTIMIT dataset is Copyright © 2001-2008 Conrad Sanderson. Distribution and research usage of this dataset is permitted under the following conditions:

1. This notice is left intact and not modified in any way.
2. The dataset is provided as is. There is no warranty as to the fitness for any particular purpose.
3. The author of the dataset is not responsible for any direct or indirect losses resulting from the use of the dataset.
4. Any publication (e.g. conference paper, journal article, technical report, book chapter, book) resulting from the usage of the VidTIMIT dataset must cite the following book:
  - C. Sanderson.  
Biometric Person Recognition: Face, Speech and Fusion.  
VDM-Verlag, 2008.  
ISBN 978-3-639-02769-3.

## Overview

The VidTIMIT dataset is comprised of video and corresponding audio recordings of 43 volunteers (19 female and 24 male), reciting short sentences. It can be useful for research on topics such as automatic lip reading, multi-view face recognition, multi-modal speech recognition and person identification/verification.

The dataset was recorded in 3 sessions, with a mean delay of 7 days between Session 1 and 2, and 6 days between Session 2 and 3. The delay between sessions allows for changes in the voice, hair style, make-up, clothing and mood (which can affect the pronunciation), thus incorporating attributes which would be present during the deployment of a verification system. Additionally, the zoom factor of the camera was randomly perturbed after each recording.

The sentences were chosen from the test section of the NTIMIT corpus<sup>1</sup>. There are ten sentences per person. The first six sentences (sorted alpha-numerically by filename) are assigned to Session 1. The next two sentences are assigned to Session 2 with the remaining two to Session 3. The first two sentences for all persons are the same, with the remaining eight generally different for each person. The mean duration of each sentence is 4.25 seconds, or approximately 106 video frames (using 25 fps).

A typical example<sup>2</sup> of the sentences used is in Table 1. There is complete correspondence of the subject IDs between VidTIMIT and NTIMIT (and hence the recited sentences).

In addition to the sentences, each person performed an extended *head rotation* sequence in each session, which allows for extraction of profile and 3D information. The sequence consists of the person moving their head to the left, right, back to the center, up, then down and finally return to center.

---

<sup>1</sup>C. Jankowski, A. Kalyanswamy, S. Basson and J. Spitz, "NTIMIT: A Phonetically Balanced, Continuous Speech Telephone Bandwidth Speech Database", Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP), Albuquerque, 1990, Vol. 1, pp. 109-112. NTIMIT can be obtained from the Linguistic Data Consortium ([www ldc upenn edu](http://www ldc upenn edu)).

<sup>2</sup>Copyright restrictions on the NTIMIT corpus prevent the list of all sentences used in VidTIMIT.

Section ID	Sentence ID	Sentence text
Session 1	sa1	She had your dark suit in greasy wash water all year
	sa2	Don't ask me to carry an oily rag like that
	si1398	Do they make class-biased decisions?
	si2028	He took his mask from his forehead and threw it, unexpectedly, across the deck
	si768	Make lid for sugar bowl the same as jar lids, omitting design disk
	sx138	The clumsy customer spilled some expensive perfume
Session 2	sx228	The viewpoint overlooked the ocean
	sx318	Please dig my potatoes up before frost
Session 3	sx408	I'd ride the subway, but I haven't enough change
	sx48	Grandmother outgrew her upbringing in petticoats

Table 1: Typical example of sentences used in the VidTIMIT database

The recording was done in a noisy office environment (mostly computer fan noise) using a broadcast quality digital video camera. For almost all the recordings the lighting setup<sup>3</sup> was as follows:

1. Standard overhead fluorescent tubes, like in most office environments. The lights were covered with A4 size white office paper in order to diffuse the light – this reduced the glare on the face and top of the head.
2. An incandescent lamp in front of the person (just below the camera). The lamp was covered with a sheet of A4 size white office paper.

The video of each person is stored as a numbered sequence of JPEG images with a resolution of  $384 \times 512$  pixels (rows  $\times$  columns). A quality setting of 90% was used during the creation of the JPEG images. The corresponding audio<sup>4</sup> is stored as a mono, 16 bit, 32 kHz WAV file.

The VidTIMIT dataset is comprised of 44 files (including this documentation) taking up about 3 Gb. Each zip archive is for one person (e.g. *felc0.tar*) and has the following internal structure:

```
subjectID / audio / sentenceID.wav
subjectID / video / sentenceID / ###
```

where *sentenceID* is the head rotation or sentence identifier (e.g. *sx396*) and *###* is a three digit<sup>5</sup> frame number (e.g. 037). Each frame is stored as a JPEG image (note that there is no *.jpg* extension). There is no audio for the head rotation sequences.

## Acknowledgements

The VidTIMIT dataset was created by Conrad Sanderson while he was a PhD student at Griffith University, Queensland, Australia (see [www.gu.edu.au](http://www.gu.edu.au)), under the supervision of Professor Kuldeep K. Paliwal.

<sup>3</sup>An exception to the lighting setup occurs for Session 3 of *mbdg0* and Session 2 of *mjar0*. For these cases the incandescent lamp was switched off.

<sup>4</sup>The audio was recorded using the camera's microphone.

<sup>5</sup>The three digit frame number limitation does not apply to the *head2* and *head3* sequences of *mrgg0*.