

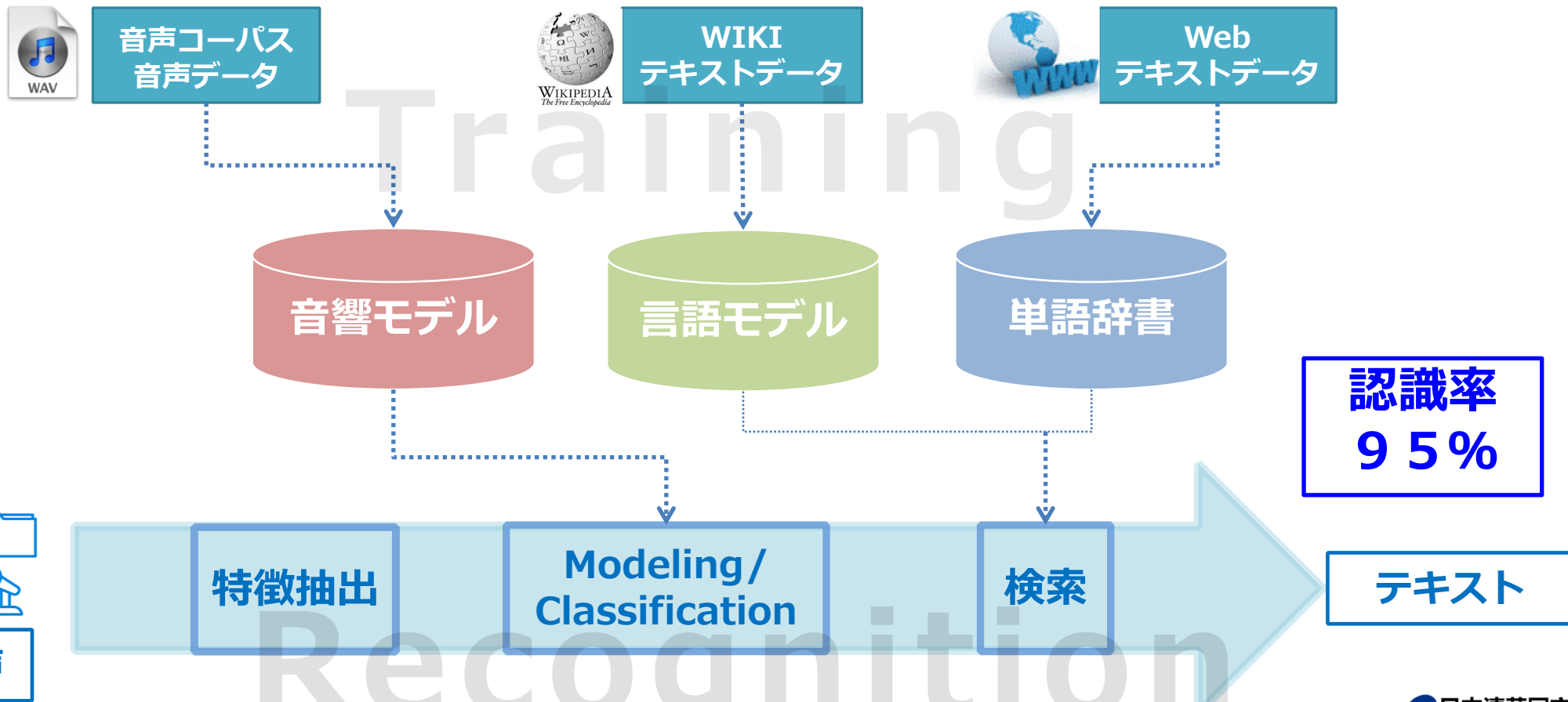
音声認識研究内容紹介



目次

- 音声認識の全体構成図
- 属性説明
- 音声認識エンジンの比較
- 知識のトランスファー（KT）技術説明
- 訓練用音声データの採取
- 運用イメージ

音声認識の全体構成図



属性説明

音響モデル	音声ドメイン	Deep Learning技術	モデルサイズ
AM	講演	○	76 MB

言語モデル	学習用データ		辞書	LMサイズ (MB)
	データ名	ドメイン		
LM	WEB 07	-	vocab-20170217	457
	WIKI	百科事典		

辞書	単語の数	選定基準
vocab	89556	WIKIに出現頻度が高い順（総数約9万まで）

デコーディング	サイズ (GB)	認識処理速度 (1秒音声データに対する)
decode	6～7	0.45秒

音声認識エンジンの比較

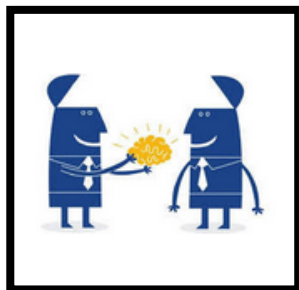
	IBM	Google	Microsoft	日本清華同方
開発元	IBM社	Google社	Microsoft社	清華同方研究新規開発
音声認識サービス運用	IBMクラウドにアップロード	Googleクラウドにアップロード	MSクラウドにアップロード	お客様の自由 クラウド運用ができる オンプレもできる
日本語音声認識	97% 個別再訓練対応しない	99% 個別再訓練対応しない	97% 個別再訓練対応しない	95.8% お客様要望より、業務向け強化再訓練が可能
A Mモデル訓練用データ資源量	数万時間	数万時間	数万時間	数百時間(*KT技法)
サービス(I/F)	サービス公開 (SDK)	サービス公開 (SDK)	サービス公開 (SDK)	SDK公開 WebService (Socket/NodeJS) お客様要望より、カスタマイズし、第三者に非公開も可能
Knowledge Transfer技法	なし	なし	なし	中国語の音声認識ニューラルネットよりKT ※

※KT(Knowledge Transfer) は中国で特許提出済、現在は中国特許庁で審査中

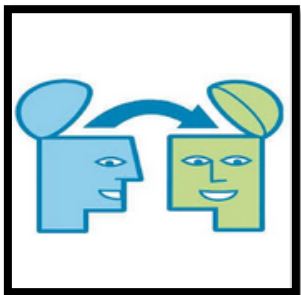
知識のトランスファー（KT）技術説明



- AIの転移学習技術です

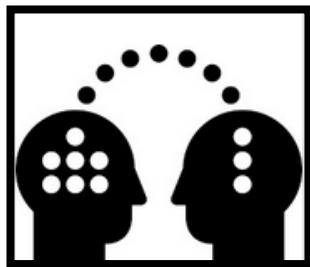


- AIはAIをトレーニングする



- 目的は学習データ準備の大変さを緩和

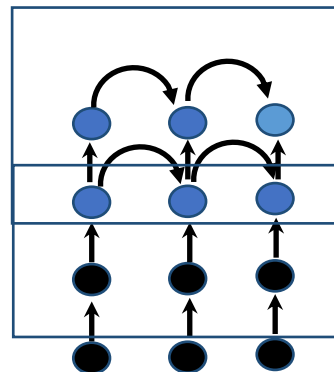
注：中国語の音声資源は5000H超、それらを強い計算環境でAM生成済み、ちなみに日本語の資源は660H



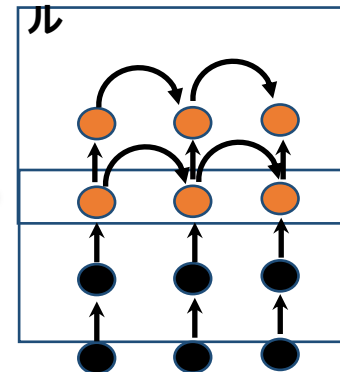
- 方言の音声認識時にモデルの作成に活用できる



中国語のAMモデル



日本語のAMモデル



講師ありのDNN（深層ニューラルネット）学習手法の下で神経細胞間と階層間の情報伝授にKTの介入で、パラメータの修正を加えもっとも有効かつ正確な調整は技術で資源データの不足と計算能力不足の

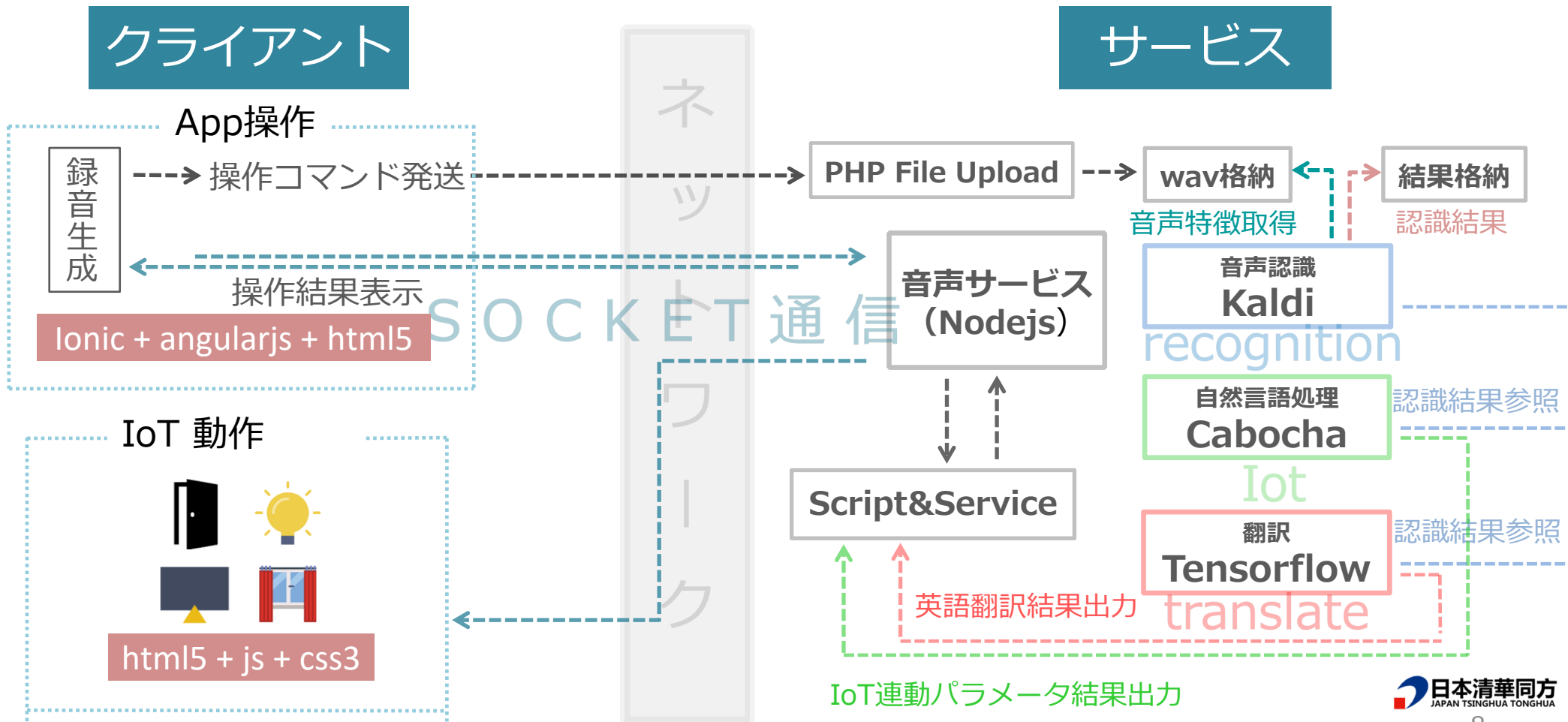
訓練用音声データの採取：区切りとラベリング※

The screenshot displays the Audacity audio editing interface. At the top, the title bar reads "Sound (A01F0055)". The main window shows a black waveform on a white background. A vertical red line is positioned at approximately 7.0 seconds. Below the waveform, a text track contains Japanese text with time markers: "えーっと) / (R xxx / xx / xx / xxxxx / xxxxx / xxxxx) という / ことで / (F エー) / (". Further down, a detailed labeling table is visible, listing time intervals and corresponding labels. To the right, a spectrogram shows the frequency content of the audio, with a time marker at 2:08.136. The bottom of the interface shows the Audacity status bar with various settings and a file list on the left.

Time Interval	Label
0008 0007.004-00008.370 L:	& (F エー)
私共は	& ワタクシドモワ
0009 0008.870-00011.195 L:	& ニュージガ
乳児が	& オンガクオ
音楽を	& ドノヨーニ
どのよ	& キーテイルカ
聞いているか	& マタ
0010 00011.401-00015.595 L:	& チョーシュニ
また	
聴取に	

※NN深層学習の性格上、大量に採取できた音声は正しくマックつけ、講師ありの訓練材料となる

運用イメージ



ご清聴ありがとうございました