# Dual Acoustic Models and Probabilistic Cross-Lingual Speaker Adaptation for Bilingual Speech Synthesis for a Monolingual Speaker

## Simple4All Consortium Submission to Blizzard-2014 Spoke Task

### I. Objective

- **Task: Bilingual Speech Synthesis for a Monolingual Speaker**

- To sythesize dual-language utterances, primarily a native language (Indian) intersperced with words from a non-native language (English)
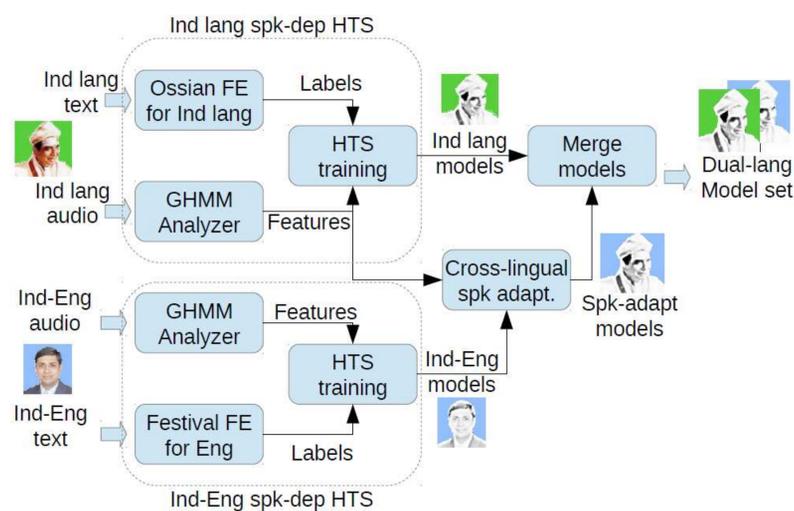
- Training data
  - Single speaker data only in Indian language (a few hundred utterances)
    - Example: "प्रसिद्ध कबीर अध्येता, पुरुषोत्तम अग्रवाल का यह शोध आलेख, उस रामानंद की खोज करता है "
  - Audio data (16kHz, 16 bits) along with text in Indian script (UTF-8)
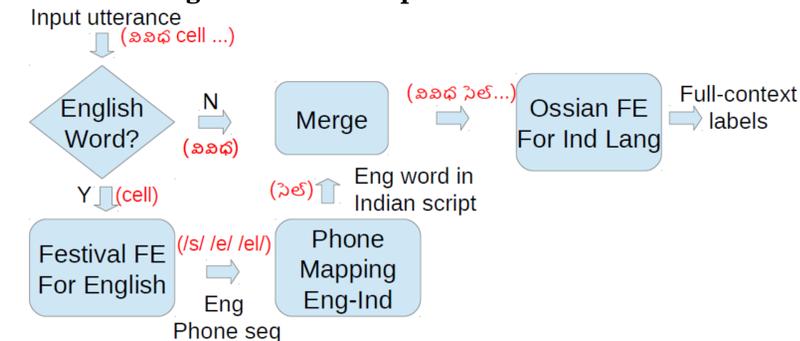- Test data
  - Example: "Under 19 cricket world cup में सोमवार को अफ़गानिस्तान ने ऑस्ट्रेलिया को हराकर, बड़ा उलटफेर किया है"
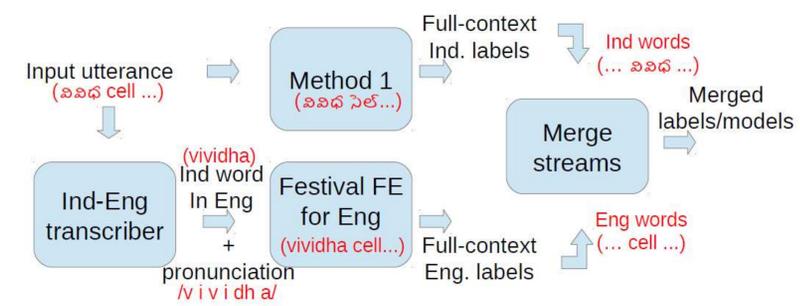
### II. Dual Language Acoustic Modeling - Overview



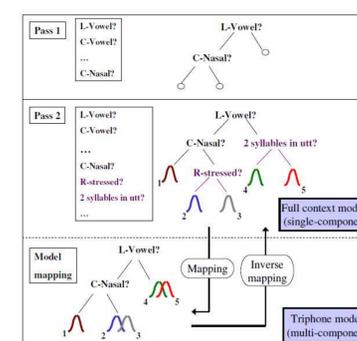### III. Label Generation

**Method-1: Eng-to-Ind transcription**



**Method-2: Dual front-end with filler words**



### IV. Acoustic Models training

- **GlottHMM based TTS models**

- **Speaker dependent models for Indian language**

- **Cross-lingual adaptation of an Indian accented English voice using Indian language data**

  - Two pass decision tree generation
  - Decode Indian language data using triphone models
  - CMLLR adaptation of the English model set
  - Merge Indian speaker-dependent and adapted English acoustic model sets

### V. Two-pass decision tree generation



- **1st pass:** Generate tree for ASR type triphones
- **2nd pass:** Extend ASR tree to TTS type full-context tree
- Train the TTS tree leaf Gaussians
- Finally pool all TTS Gaussians under each ASR tree leaf to form Gaussian mixture models

### VI. Demo Page

**http://research.ics.aalto.fi/speech/demos/COIN_blizzard14/**