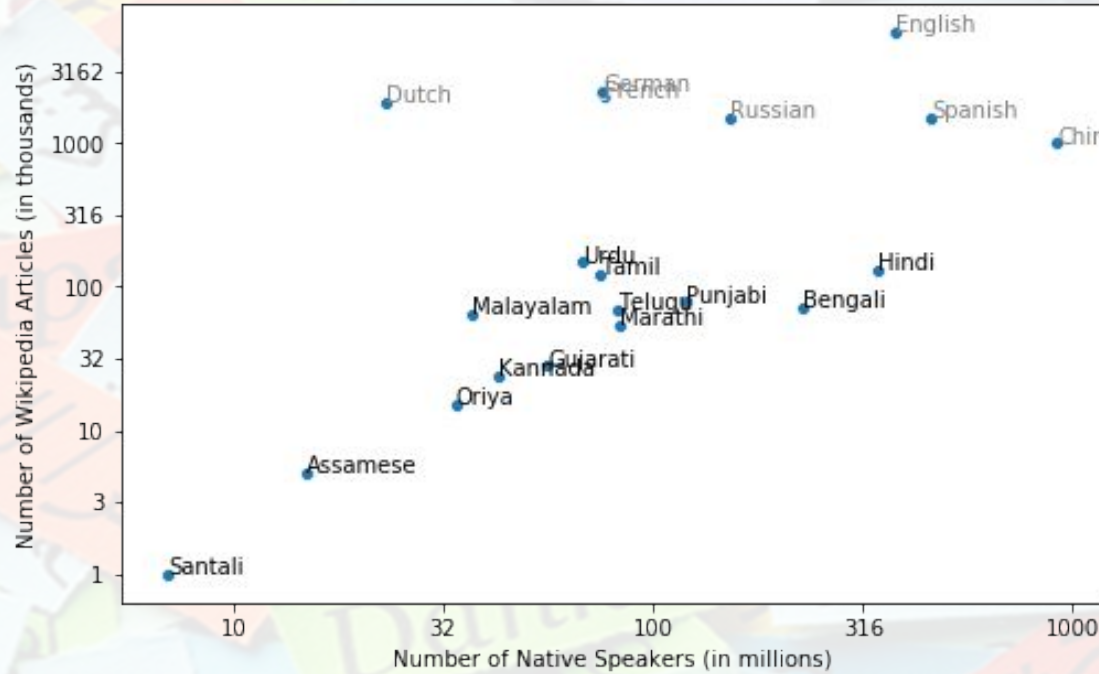# Unsung Challenges of Building and Deploying Language Technologies for Low Resource Language Communities

Pratik Joshi, C. Barnes, Sebastin Santy, S. Khanuja, S. Shah, A. Srinivasan, S. Bhattamishra, S. Sitaram, K. Bali, M. Choudhury

**Microsoft Research India**

# Resource Disparity

# Discussion Areas

- Information Exchange
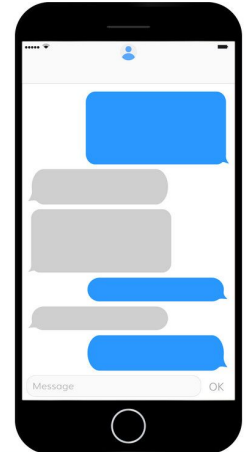  - Access
  - Generation
- Interface
- Deployment and Impact

Microsoft

# Information Exchange

Microsoft

# Access of Information

Making existing information accessible to people:

- Must know : Disaster/Banks
- Should know : Rights/Duties/Health
- Can know : Agriculture/Financial Management
- Want to know: Forms of Entertainment

Microsoft

# Access of Information

Making more information digitally available to more people:

- Make effective translation tools

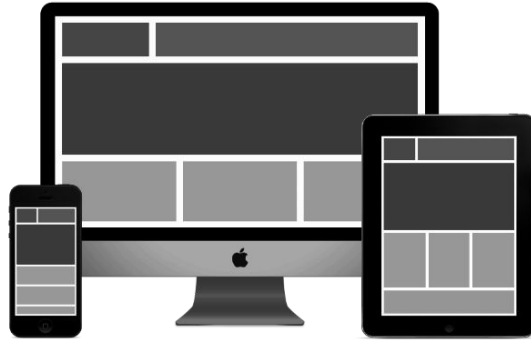Making NLP models more accessible to low resource languages:

- Difficulty in adopting NLP tools to low resource languages
- Use of Machine Translation tools

Microsoft

# Generation of Information

Making non-tangible information accessible to people:

- Digitization of Documents
- Crowd-sourcing

Microsoft

Interface

Microsoft

# Issues with Text-based Interface

- Thies et.al., 2015 reports that text-based interfaces are:
  - Redundant for illiterate users,
  - Severely error-prone for novice, literate users.
- Several languages don't have unique keyboard/fonts, some lack script overall (Boyera, 2007).

Microsoft

# Alternate Modalities - Speech



- CGNet-Swara (Mudliar et al., 2013)
  - Phone-based IVR System
  - Citizen-run journalism portal for educating illiterate users
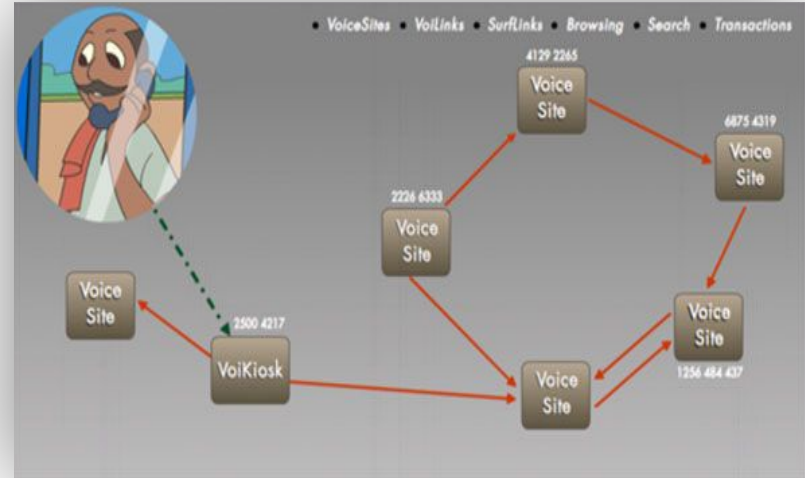
# Alternate Modalities - Speech


avaajotalo
Avaaj Otalo - voice-based social software

- Avaaj Otalo (Patel et al., 2010)
  - Simple Phone Call System
  - Browse or ask questions on agricultural topics



Microsoft

# Alternate Modalities - Speech



- Spoken Web (Kumar et al., 2010)
  - Voice interface application
  - Enables creation of voice sites analogous to websites
  - Voice sites accessible through application

# Alternate Modalities - Voice+Graphic Based



- VideoKheti (Cuendet et al., 2013)
  - Multi-modal mobile interface with large buttons, graphics, voice input
  - Provides agricultural videos in desired dialect

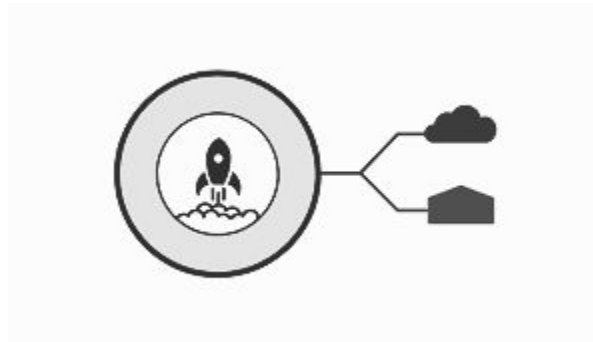# Alternate Modalities - Voice+Graphic Based


Google bolo

- Bolo
  - Speech-based reading tutor application
  - Helps improve reading skills for children



स्कूल में आज मेरा पहला दिन है।
माँ मेरा हाथ पकड़े हुए हैं और मेरे साथ चल रही हैं।

"मैं अब बड़ी हो गई हूँ," मैं कहती हूँ।

"चलो...चलो!"

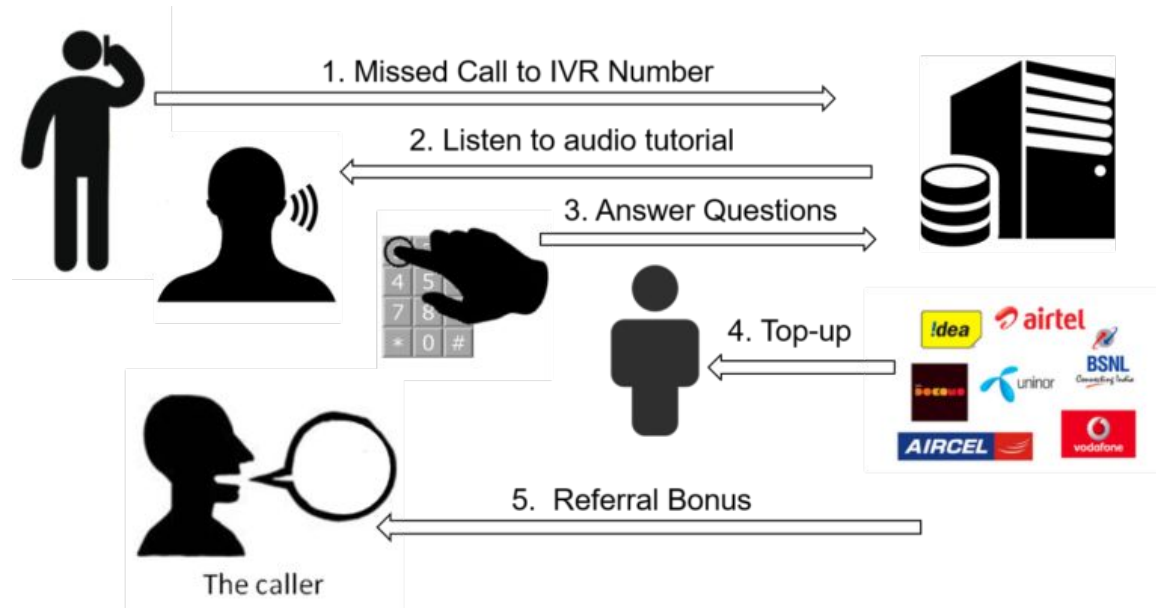माँ ने मेरा हाथ कसकर पकड़ा हुआ है।

Microsoft

# Challenges to Voice Interfaces

- Difficult to build robust ASR systems for these languages due to:
  - Lack of data
  - Dialect variations
- Qiao et al., 2010 attempts to solve this via SALAAM ASR:
  - Cross-lingual phoneme mapping between high-to-low resource language
  - Limited vocabulary, but cost-effective
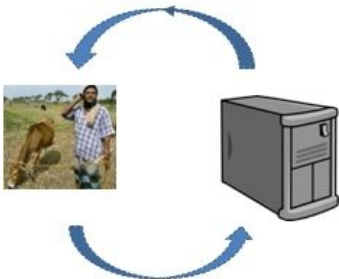
Microsoft

# Deployment and Impact

# Learn2Earn

# Learn2Earn

- Tech-enabled Information dissemination system
- Deployment
  - Originally centred on informing farmers about their rights (Indian Forest Rights Act)
  - Advertisement seeding on existing IVR channel
  - Financial incentives to use platform, invite friends
- Impact
  - Platform spread from initial 17 users to over 17,000

# Mobile Vaani



**1. Speak**

Users speak and listen to contributions over our intelligent IVR platform

**2. Moderate + Share**

Content moderated locally and centrally, then published on IVR, web
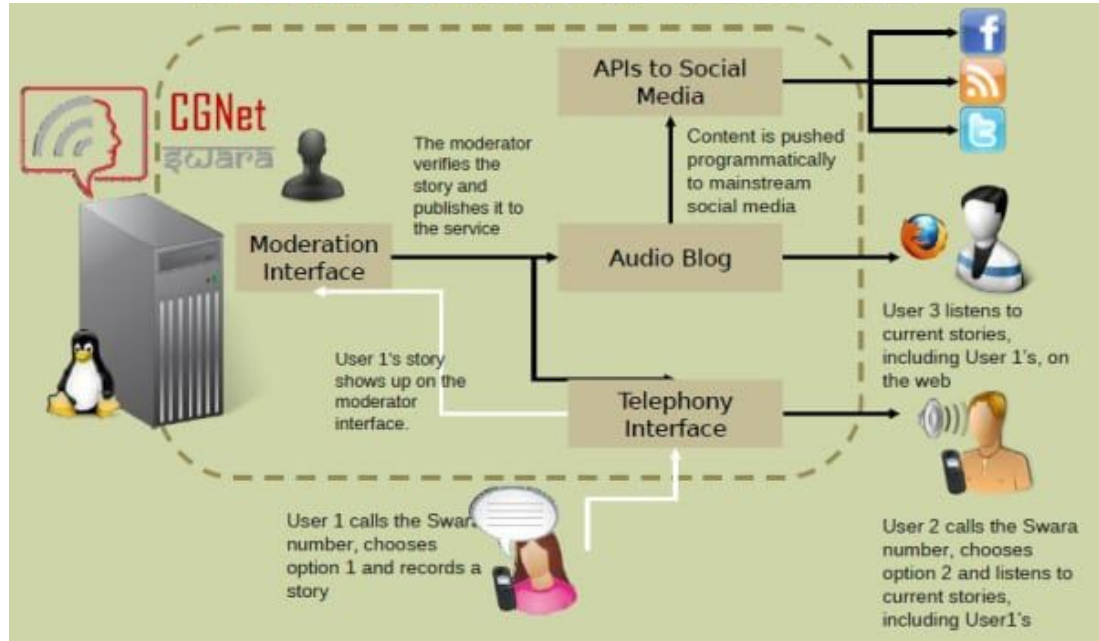
**3. Connect to stakeholders**

Inputs connected to government (local + other), NGO partners, social enterprise partners

# Mobile Vaani

- Voice-based medium for digital communication, information sharing
- Deployment
  - First employees recruited did not originate but worked closely with community
  - Informed end-users about design and use of technology
  - Re-launch involved compensation plans for trained volunteers
- Impact
  - Initial success followed by lack of growth due to unrealistic expectations in minds of users.
  - Re-launch resulted in better use and adaptation.
  - Gained popularity during teachers' strike in Jharkhand

Microsoft

# CGNet Swara

# CGNet Swara

- Citizen-run journalism portal for educating illiterate users
- Deployment:
  - ~50 workshops
  - Trained >2000 members of various communities
  - Outreach activities : financial incentives provided
- Impact:
  - Incentives empowered users to deliver the technology to desired areas
  - Outreach activities increased spread of awareness via word-of-mouth

# Conclusion

Microsoft

- NLP research breaks down communication and information barriers
- Increased investment needed for enabling low-resource communities with language technologies
- Hope that increased exposure to problem sparks more discussion and effort in this area

Microsoft