

# Spoken Language Understanding on the Edge

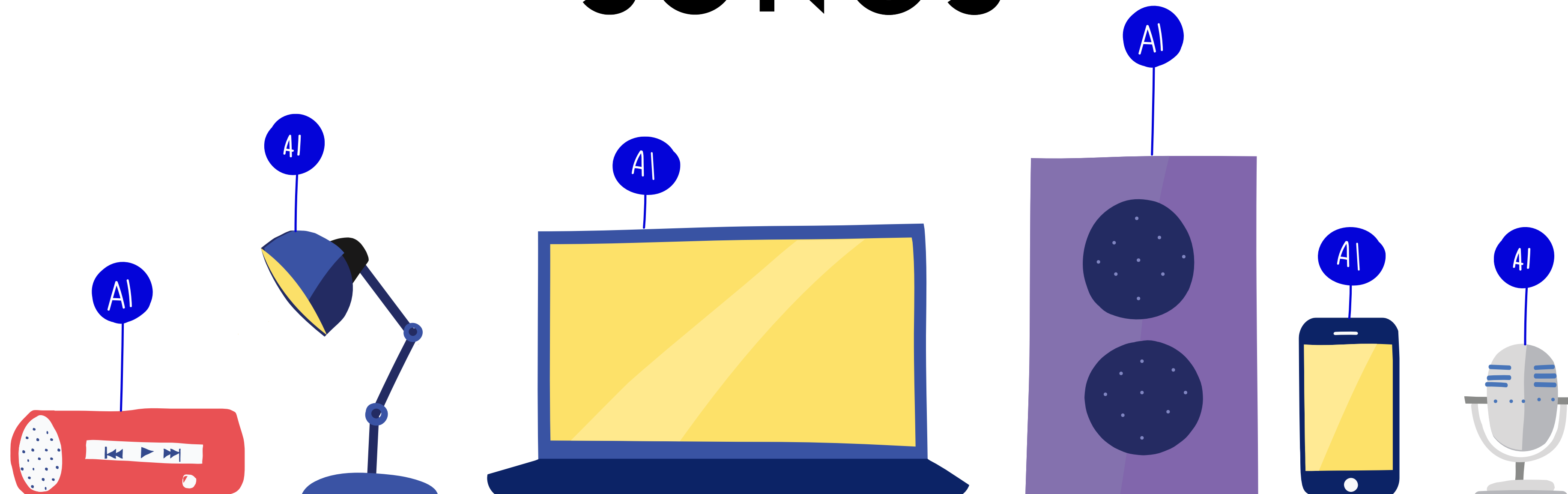
*Alaa Saade, Alice Coucke, Alexandre Caulier, Joseph Dureau, Adrien Ball, Théodore Bluche, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, Mael Primet*  
*Snips, Paris*

**EMC2 Workshop @ Neurips 2019**

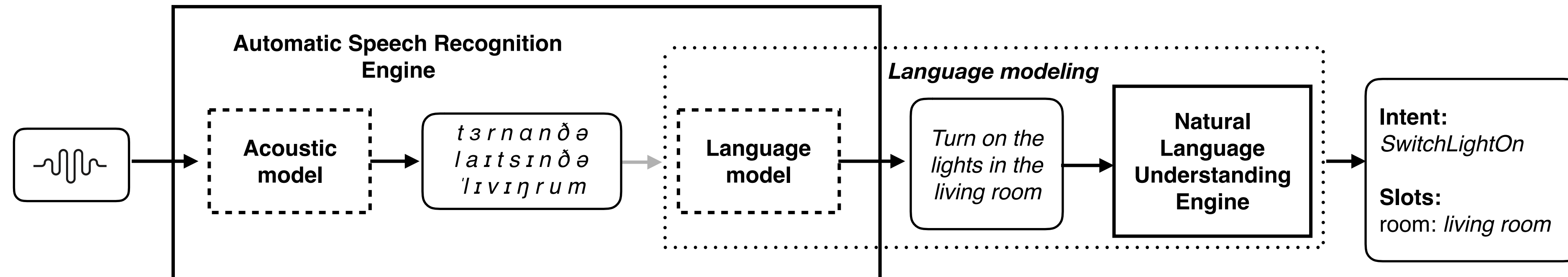
November 13

Alexandre Caulier

# SONOS



# Spoken language understanding system

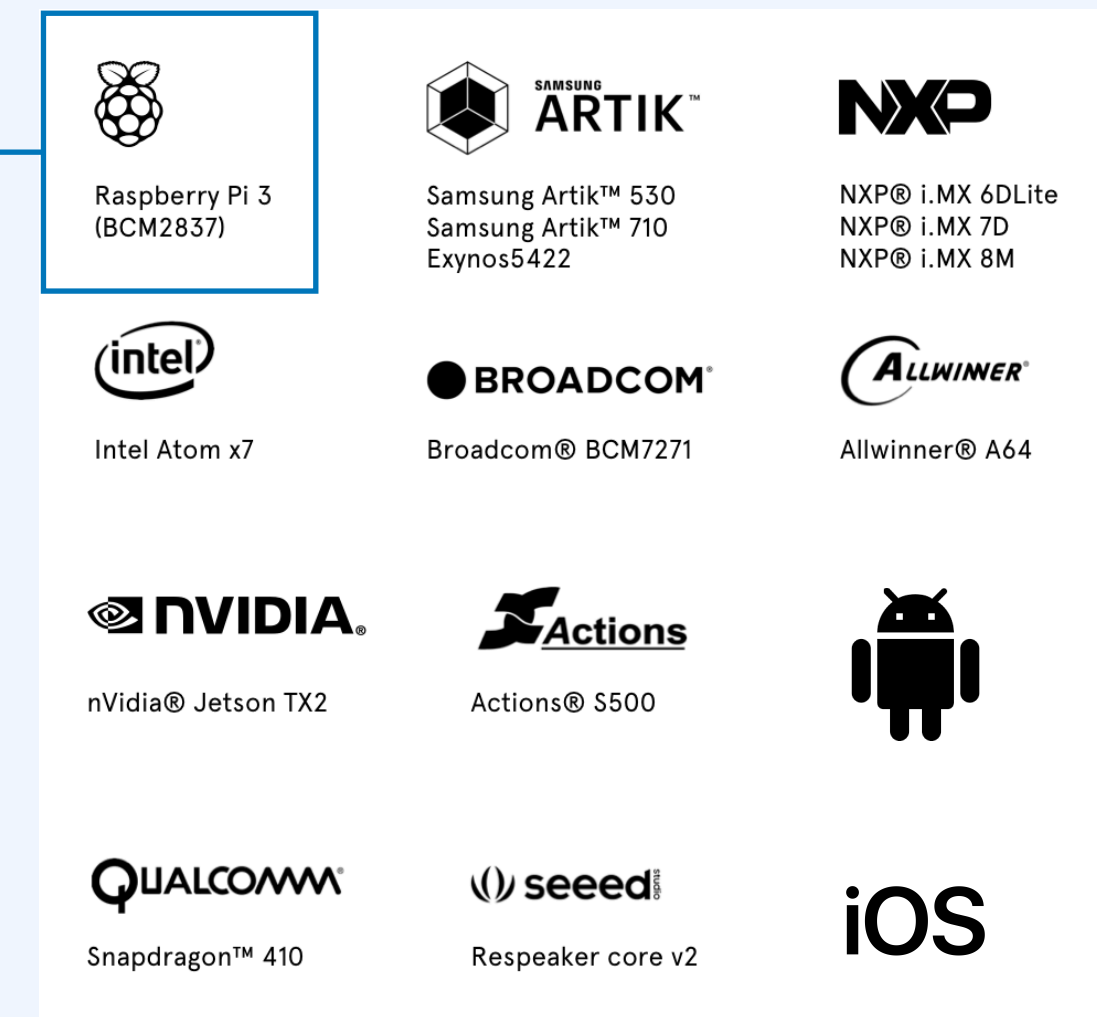


## Features

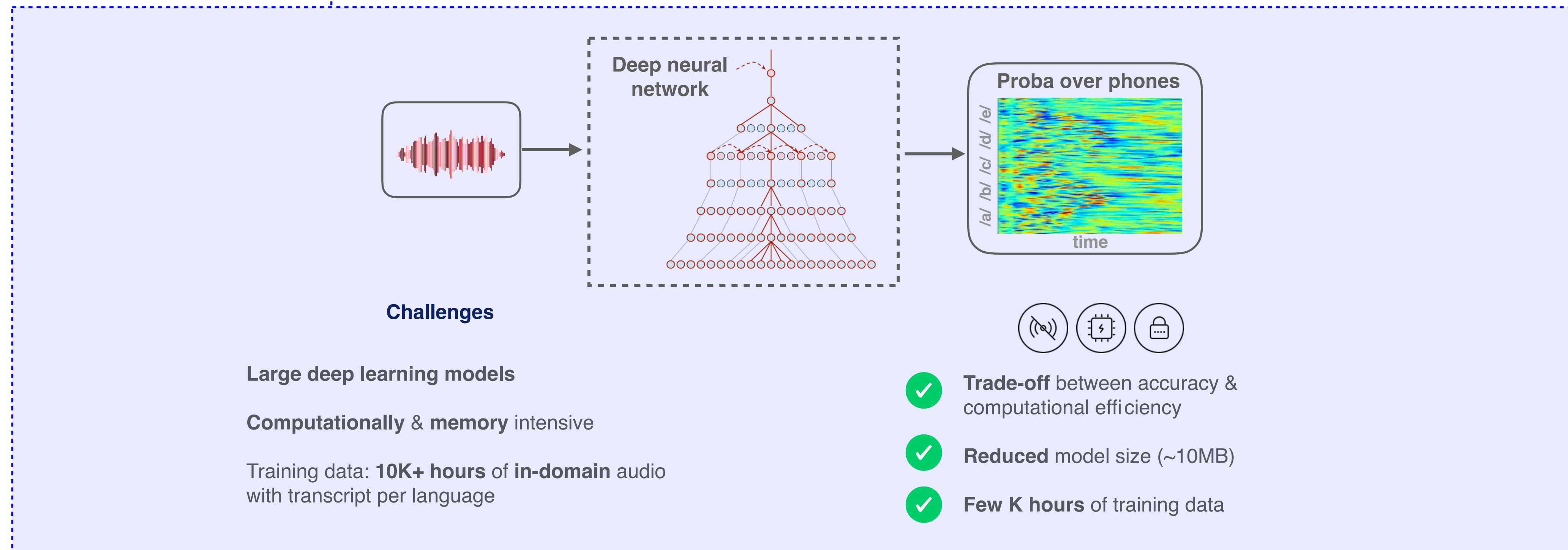
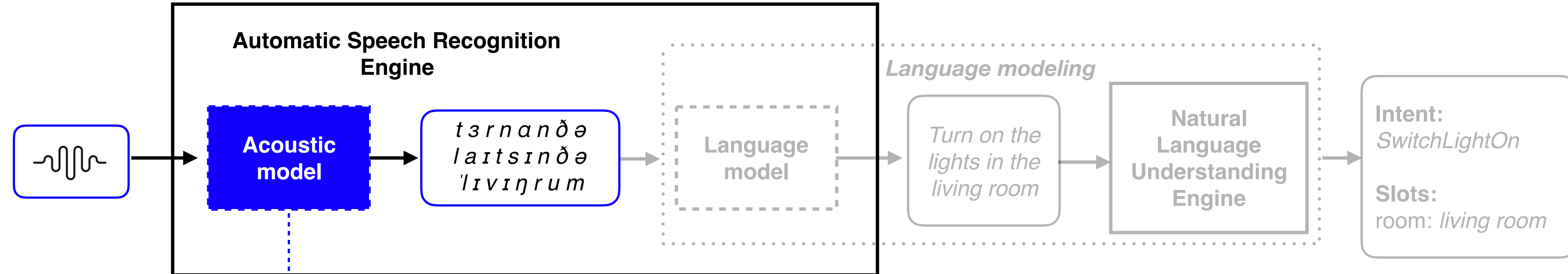
- **Cloud independent** - no remote processing
- **Private by Design** - no user data can be collected
- **Accurate** - on-par with cloud-based solutions

## Tested and certified to run on

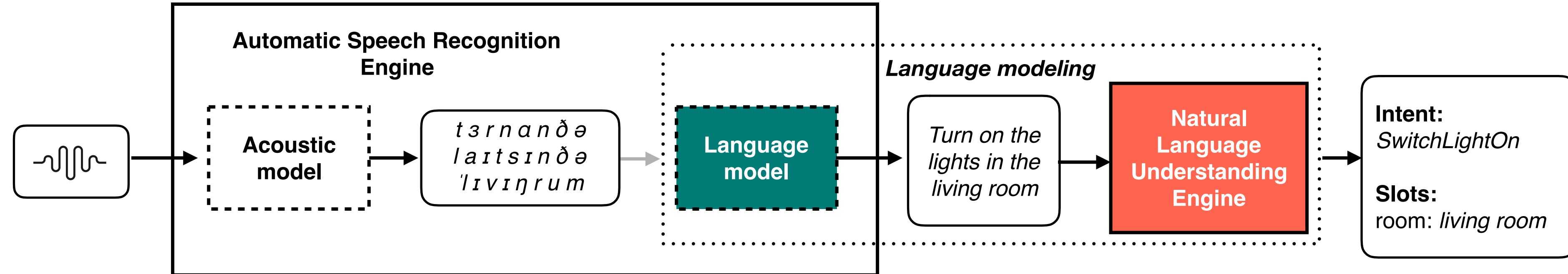
1GB RAM  
1.4GHz CPU



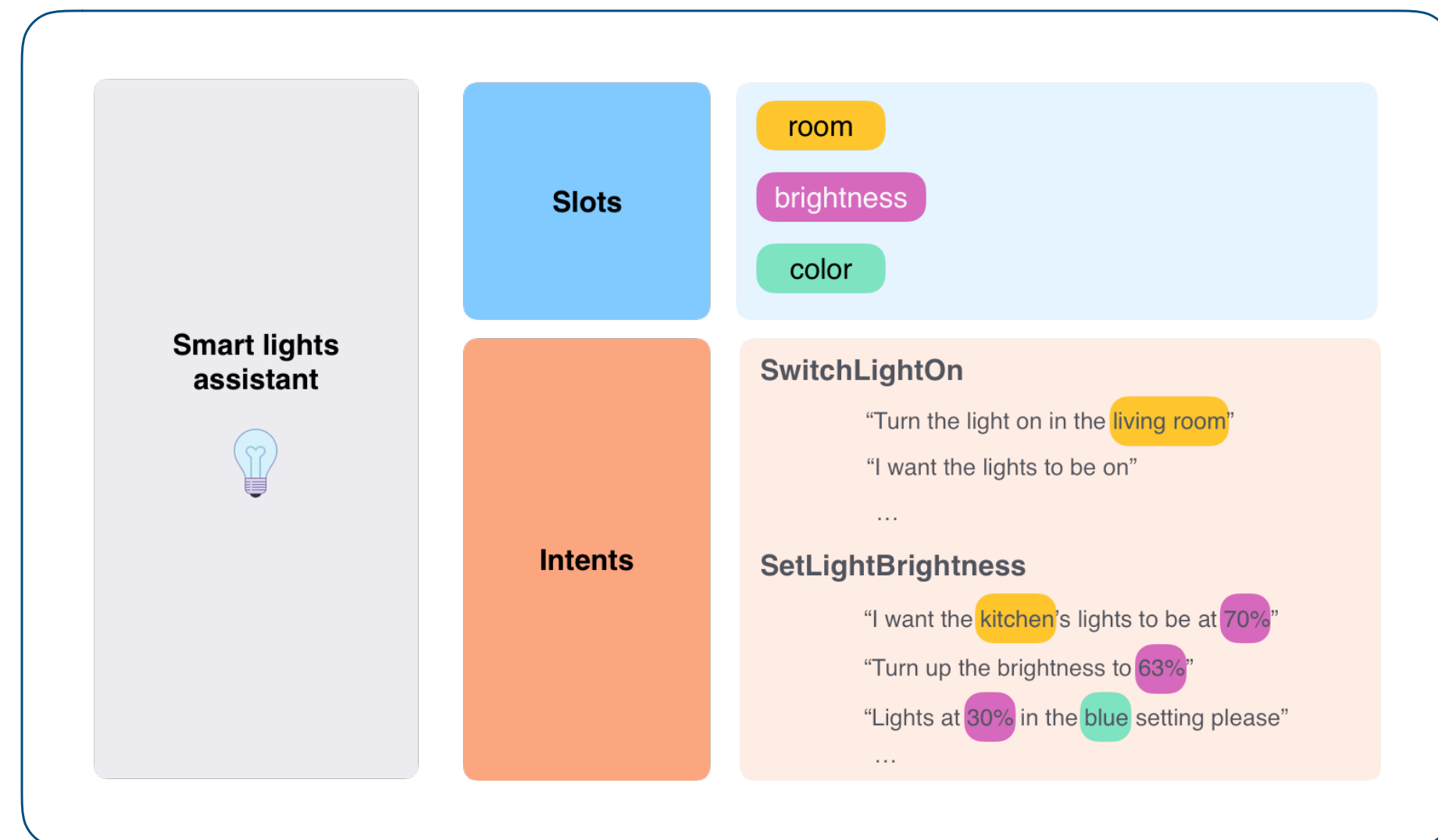
# Acoustic modeling



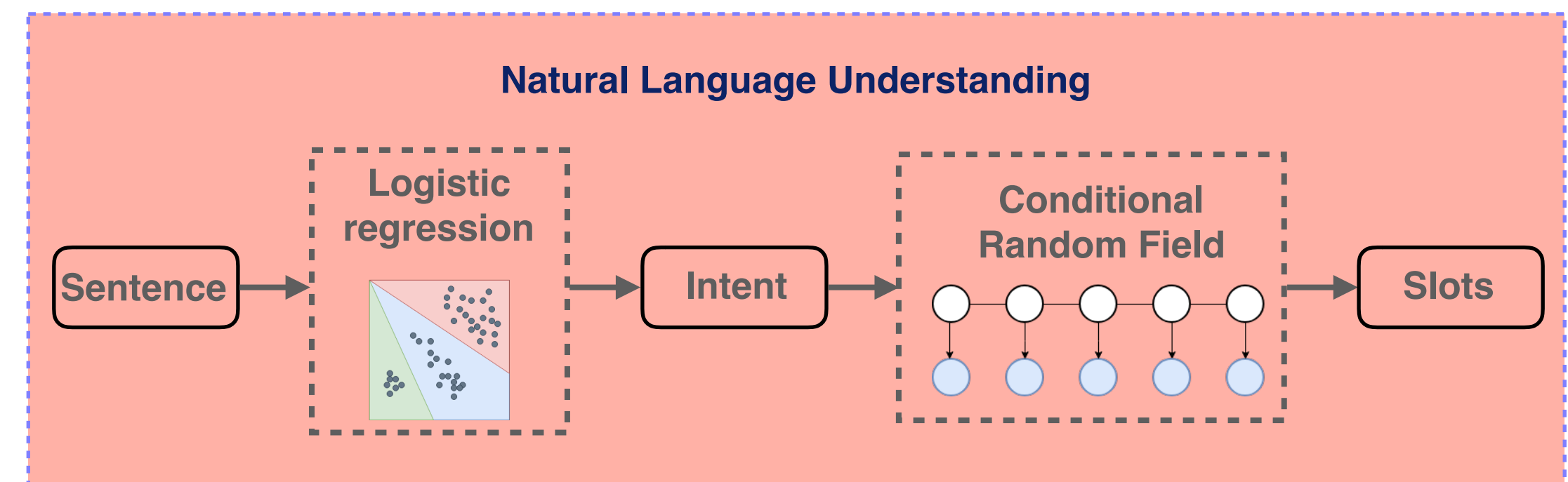
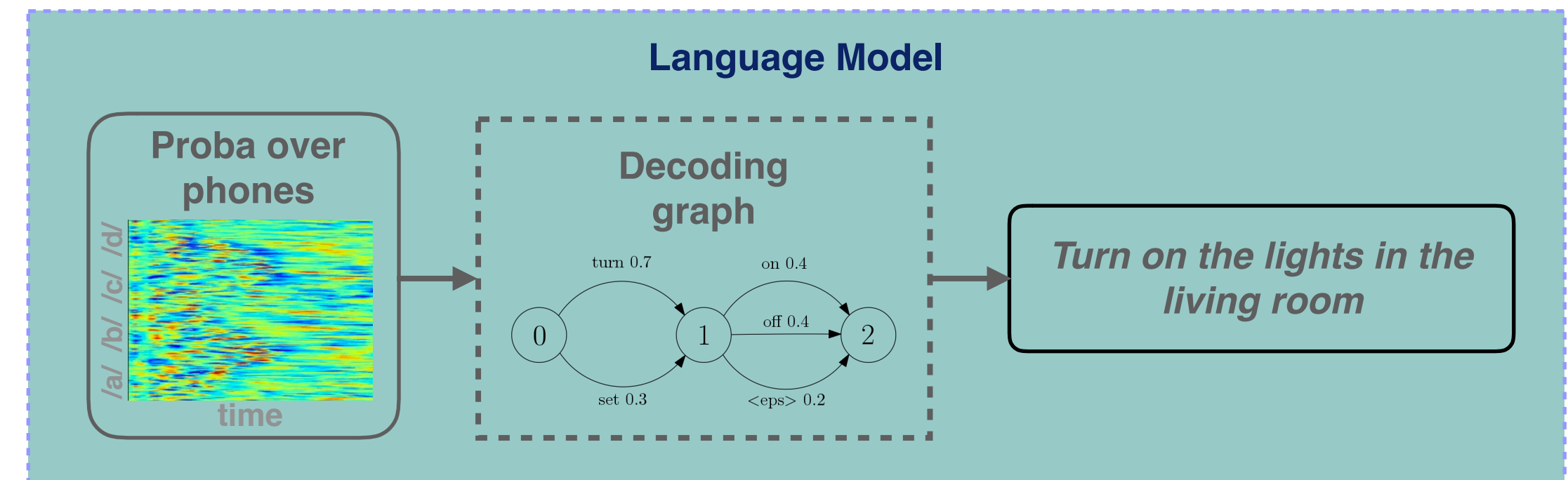
# Assistant Contextualization



**Approach :** LM and NLU are **consistent** and **contextualized**



- ✓ **Lightweight** models
- ✓ **Out of vocabulary** management
- ✓ **On-device** personalization



# Benchmarks - Datasets Open Sourcing

Experimental setting



## Datasets

Audio utterances with transcripts & supervision

Recorded in close and far-field



Smart Lights Assistant

*1.8K utterances*

*400 word pronunciations*



Music Assistant

*3K utterances*

*178K word pronunciations*

## Method



**snips**

Specialized for  &   
<100MB, real time on a  
Raspberry Pi 3



**Google**

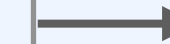
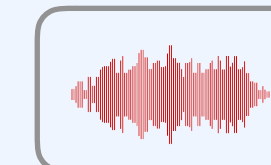
Google Speech-to-Text cloud services

*One-size-fits-all engine*

## Metrics

End-to-end score

*% of perfectly parsed queries*



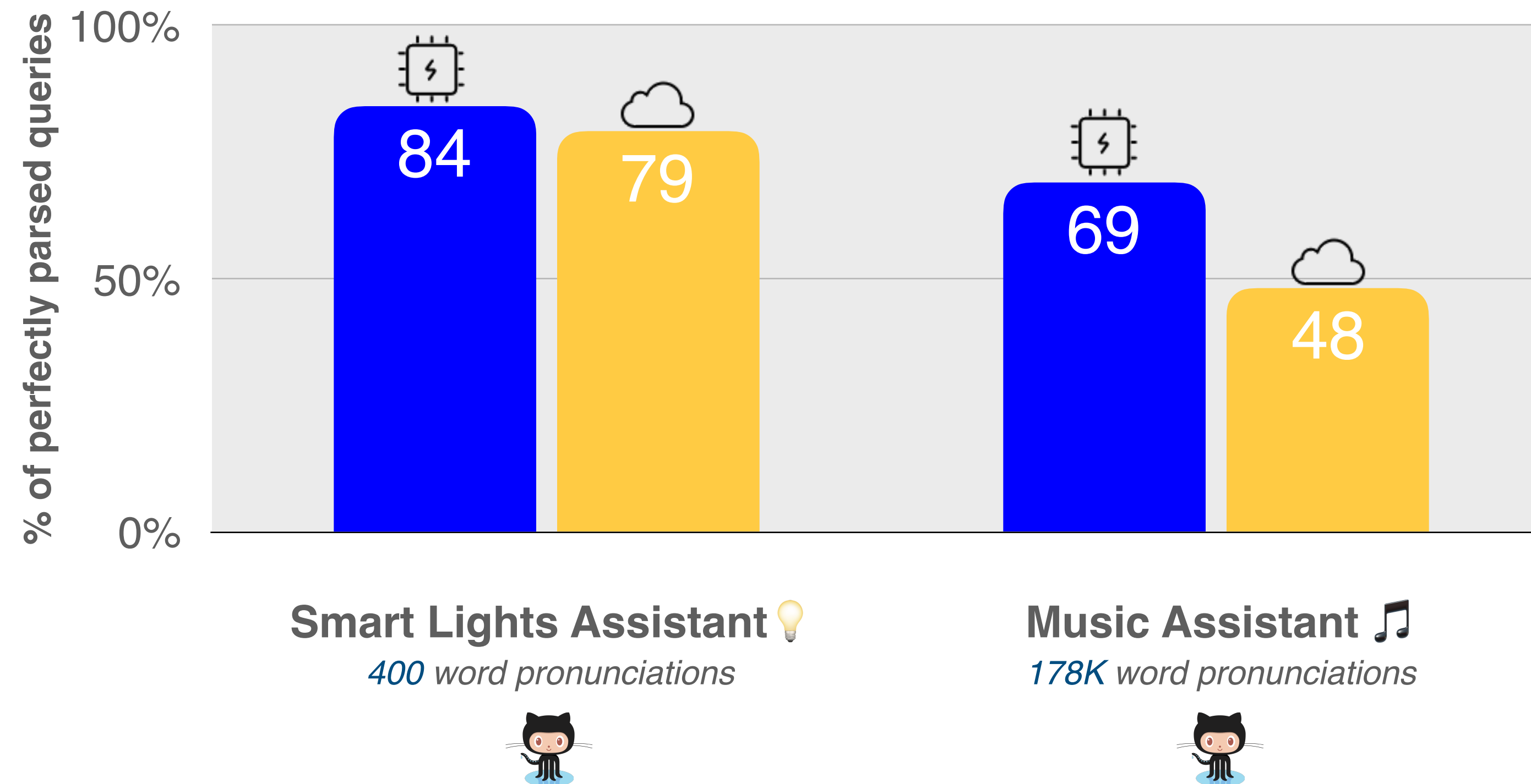
**Intent:**  
SwitchLightOn

**Slots:**  
room: living room



# Benchmarks

End-to-End performance



 **snips**

Contextualized for 💡 & 🎵  
<100MB, real time on a Raspberry Pi 3

 **Google**

STT cloud service  
One-size-fits-all engine



	Tier 1 Artists 1-1k	Tier 2 Artists 4.5k-5.5k	Tier 3 Artists 9k-10k
Snips	71 %	68 %	67 %
Google	69 %	38 %	37 %

## Questions ?