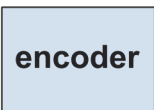
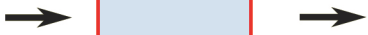
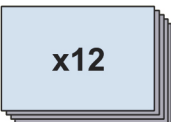




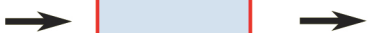
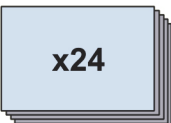
w2v



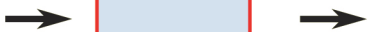
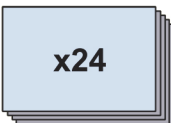
vqw2v



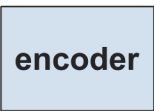
w2v2



XLSR



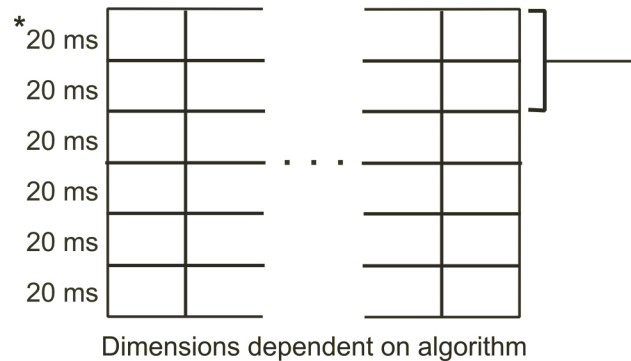
DeCoAR



MFCC



Model selection



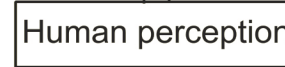
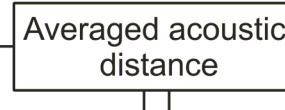
Forced alignment

| | Start | End |
|--------|--------|--------|
| PLEASE | 0.4800 | 0.8800 |
| CALL | 0.9802 | 1.1888 |
| STELLA | 1.2890 | 1.4890 |
| ASK | 2.1670 | 2.4868 |
| HER | 2.4868 | 2.5965 |
| TO | 2.5965 | ... |
| ... | | |

* : 10 ms for w2v, vqw2v, DeCoAR, and MFCC

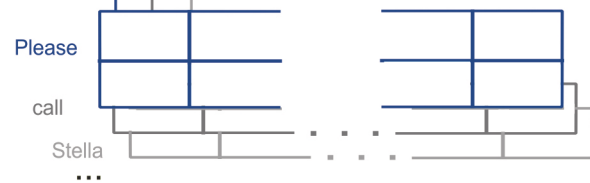
Feature extraction

Speaker 1



Pearson's correlation

Speaker 2



Dynamic Time Warping