

PROBLEM

- The success of speech-based technology relies on an understanding of speech communication.
- Common approaches:
 - linguistic models (phonemes, syllables, words, sentences).
 - speech production models (source-filter theory, articulatory phonetics).
 - auditory models (cochlear filtering, equal loudness curves, hair cell transduction).
- Needed: a model of human speech communication derived from first principles, i.e., information theory.

CONTRIBUTIONS

- An upper bound on the information rate of speech of the order of 100 b/s.
- A model of speech communication that does not rely on prior knowledge of the transmitter (vocal tract), the receiver (auditory system), or language.
- The model only relies on having recordings of multiple talkers saying the same utterance.

PREVIOUS WORK

- The probabilities of English phonemes give a lexical information rate of approximately 50 b/s [Fano, 1950].
- Variables related to talker identification, emotional state, and prosody, vary relatively slowly in time and contribute little to the overall information rate.
- The bandwidth of the human auditory system and the SNR required for perfect intelligibility give a channel capacity of approximately 20000 b/s [Flanagan, 1972].
- Why are information rates based on acoustics orders of magnitude larger than those based on linguistics?
- Hypothesized that talker variability is a type of 'production noise' that limits information transfer [Kleijn *et al.*, 2015].

HUMAN COMMUNICATION

- A talker randomly selects a message, $\{M_t\}$, e.g., a phoneme, word, or neural state, where t is the time index.
- The talker encodes the message into an acoustic speech signal, $\{S_t\}$, according to a conditional probability distribution:

$$P_{\{S_t\}|\{M_t\}}(\{S_t\}|\{M_t\}).$$

- Define a *chorus* as a set of J speech signals where each signal contains the same message:

$$\{Z_M\} = \{\{S_{M,t}\}^{(1)}, \{S_{M,t}\}^{(2)}, \dots, \{S_{M,t}\}^{(J)}\},$$

- Define a chorus-based estimate of the message as $\{\tilde{M}_t\} = f(\{Z_M\}) + \{N_t\}$ where $f(\cdot)$ is a deterministic function, $\{N_t\}$ is regularization noise.

INFORMATION BOTTLENECK

- A natural objective for the estimator $f(\cdot)$ is that it minimizes the information bottleneck:

$$f^* = \arg \min_f I(\{Z_M\}; \{\tilde{M}_t\}) - \beta I(\{S_{M,t}\}; \{\tilde{M}_t\}), \quad (1)$$

- $I(\{Z_M\}; \{\tilde{M}_t\})$ is the mutual information rate between the chorus and the chorus message estimate.
- $I(\{S_{M,t}\}; \{\tilde{M}_t\})$ is the mutual information rate between the speech and the message estimate.
- β is a Lagrange multiplier.
- Main idea: the bottleneck discards features from the chorus that aren't consistent across talkers.

COMPARISON OF ESTIMATORS

- Confine the message estimator to be of the form $f(\{Z_M\}) = \frac{1}{J} \sum_j g(\{S_{M,t}\}^{(j)})$.
- Consider $g(\cdot)$ as the identity function, STFT, spectrogram, log-spectrogram, and auditory spectrogram. Of these candidate functions, the auditory spectrogram gives the lowest bottleneck. Hence, we represent speech as a sequence of auditory-spectra: $\{X_t\} = g(\{S_{M,t}\})$.
- Suggests that the structure of speech might be adapted to the coding capability of the mammalian auditory system.

THE INFORMATION RATE

- Describe speech communication by

$$\{X_t\} = \{\tilde{M}_t\} + \{\tilde{P}_t\}, \quad (2)$$

X_t is the auditory-spectra of the speech, \tilde{M}_t the estimated message, \tilde{P}_t is Gaussian production noise.

- The mutual information rate is

$$I(\{X_t\}; \{\tilde{M}_t\}) = \lim_{k \rightarrow \infty} \frac{1}{k} I(\mathbf{X}^k; \tilde{\mathbf{M}}^k), \quad (3)$$

where \mathbf{X}^k and $\tilde{\mathbf{M}}^k$ are formed by stacking k consecutive spectra. This means that time and frequency dependencies are accounted for.

- If time-dependencies span no more than L samples, i.e., X_t is independent to X_{t+L} , then the mutual information rate reduces to

$$I(\{X_t\}; \{\tilde{M}_t\}) = \frac{1}{L} (h(\mathbf{X}^L) - h(\tilde{\mathbf{P}}^L)), \quad (4)$$

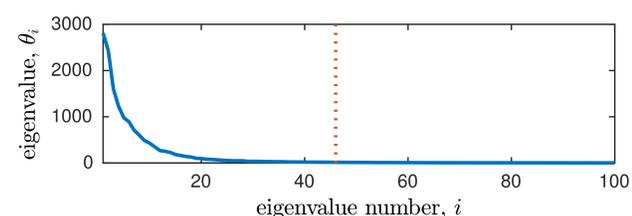
- Use Principal Component Analysis to find a subspace containing the message estimate (called the *message articulation space*).
- Reduce dimensionality by projecting stacked spectra onto the message articulation space.
- The capacity of the speech communication channel is

$$C = \frac{F}{2L} \sum_{v=1}^V \log_2 \frac{\lambda_v}{\psi_v}, \quad (5)$$

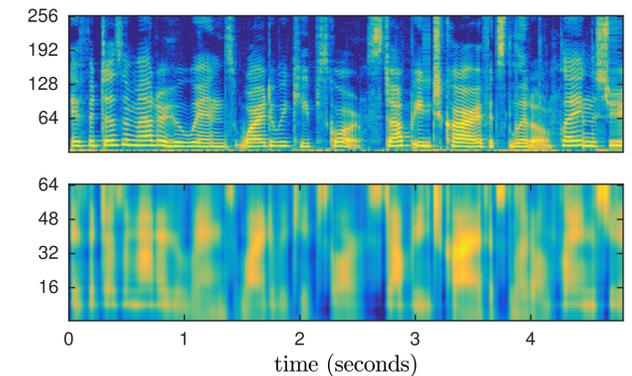
where F is the frame rate and λ_v and ψ_v are the eigenvalues of the covariance matrices of the speech and the production noise after projecting onto the message articulation space.

MESSAGE ARTICULATION SPACE

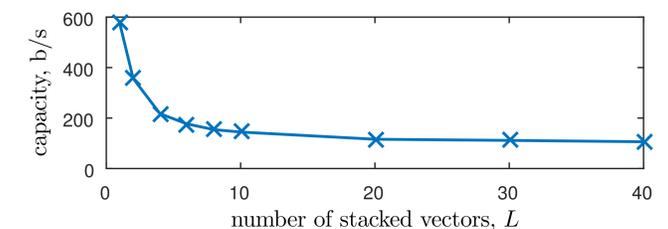
- The message estimate lies within an approximately 50-dimensional subspace:



PROJECTED SPEECH SPECTRA



CHANNEL CAPACITY



RESULTS

- The capacity of the speech communication channel is of the order of 100 b/s, which is comparable to the lexical information rate.
- Time-dependencies are negligible for $L > 10$. This corresponds to a duration of 80 ms, which is consistent with the average duration of a phoneme.

DISCUSSION

- Phonemes and words are not closed under addition. In reality the message articulation space is a manifold, not a subspace. This would lead to a lower information rate.
- In theory, given enough data, our model can also account for phoneme and word dependencies.
- Of all the representations of speech we tried, the auditory-spectra gave the lowest bottleneck, but there could be another representation that achieves a lower bottleneck. Such a representation could be found by solving (1).