# **Stream Generation Pipeline**



## International Phoneme Data

## Phonological Feature Matrix

Characterize phonemes using binary matrix of phonological features

## **Phonemes Register**

### From Phonemes to Syllables

- · Define syllable type
- Filter syllables based on unigram probability of occurence in linguistic corpus

## **Syllables Register**

## From Syllables to Pseudowords

- Define number of syllables per pseudoword
- · Generate words based on language-specific phonotactics
- Filter pseudowords based on bigram and trigram probability of occurence and positional probabilities in linguistic corpus

## **Pseudowords Register**

#### From Pseudowords to Lexicons

- Compute matrix of position-wise feature repetitions between pseudowords
- Define number of words per lxicon
- Define number of compatible lexicons
- Generate lexcons with minimal feature overlap across words and no syllable repetitions

#### Lexicons

#### From Lexicons to Streams

Generate streams with controlled transotional probabilities between syllables

## **Streams**

# **ARC Data Model**



## International Phoneme Data

```
# phonemes with phonological feature annotations
phonemes = load_phonemes()
```

## **Phonemes**

# Syllables

```
# Print output example
In [0]: print(syllables)
Out [0]: c?:|cq:|cw:|cm:|cj:|cw:|cu:|... (2294 elements total)
```

```
words = make_words(syllables, num_syllables=3,
bigram_control=True,
    trigram_control=True, phonotactic_control=True,
    positional_control=True)
```

## **Pseudowords**

## Lexicons

streams = make\_steams(lexicons, max\_rhythmicity\_index=0.01)

## **Streams**

# -









