

Hadith Segmentation Using N-grams

Introduction to NLP (EECE 634) Project Report

Ahmad Mustapha¹, Joseph Sabbagh²

¹ Electric and Computer Engineering, AUB

²Computer and Communication Engineering, AUB

Context

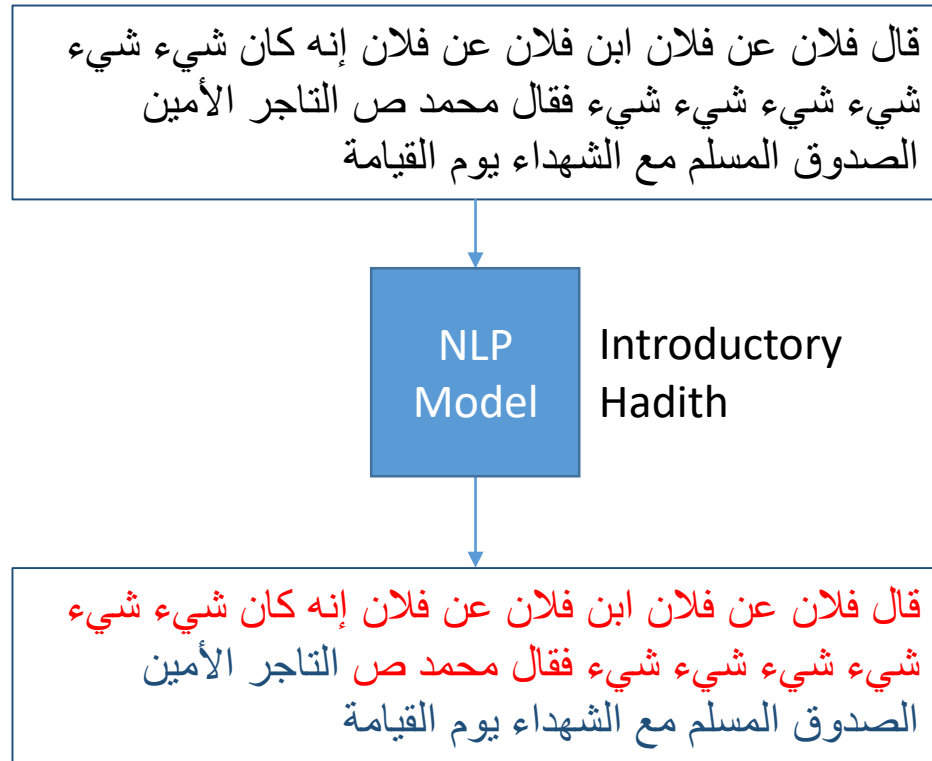
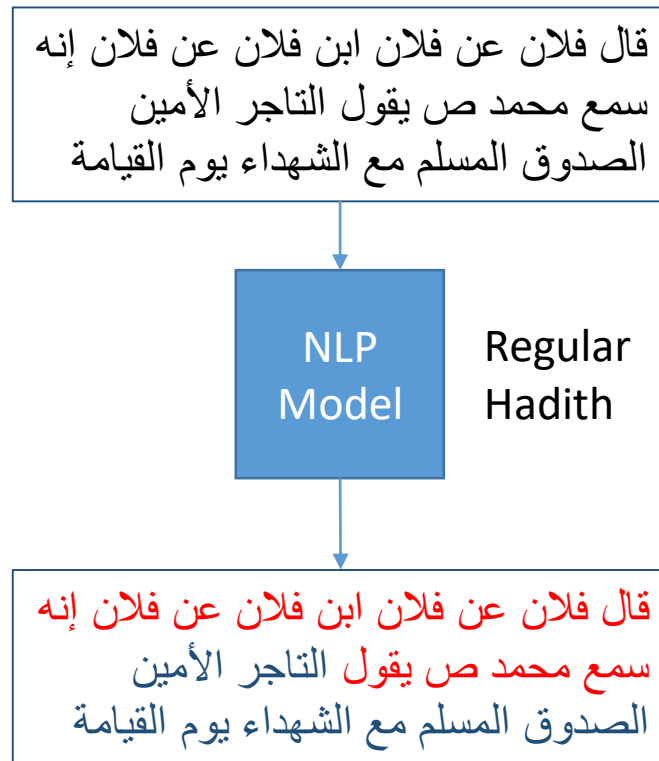
المتن

وَحَدَّثَنِي ابْنُ نُفَيْرٍ، حَدَّثَنَا أَبِي، حَدَّثَنَا حَنْظَلَةُ، قَالَ سَمِعْتُ
عِكْرَمَةَ بْنَ خَالِدٍ، يُحَدِّثُ طَاوُسًا أَنَّ رَجُلًا، قَالَ لِعَبْدِ اللَّهِ بْنِ
عُمَرَ أَلَا تَغْزُو فَقَالَ إِنِّي سَمِعْتُ رَسُولَ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ
وَسَلَّمَ يَقُولُ " إِنَّ الْإِسْلَامَ بُنِيَ عَلَى خَمْسٍ شَهَادَةِ أَنْ لَا إِلَهَ إِلَّا
اللَّهُ وَإِقَامَ الصَّلَاةِ وَإِيتَاءَ الزَّكَاةِ وَصِيَامِ رَمَضَانَ وَحَجِّ الْبَيْتِ "



الإسناد

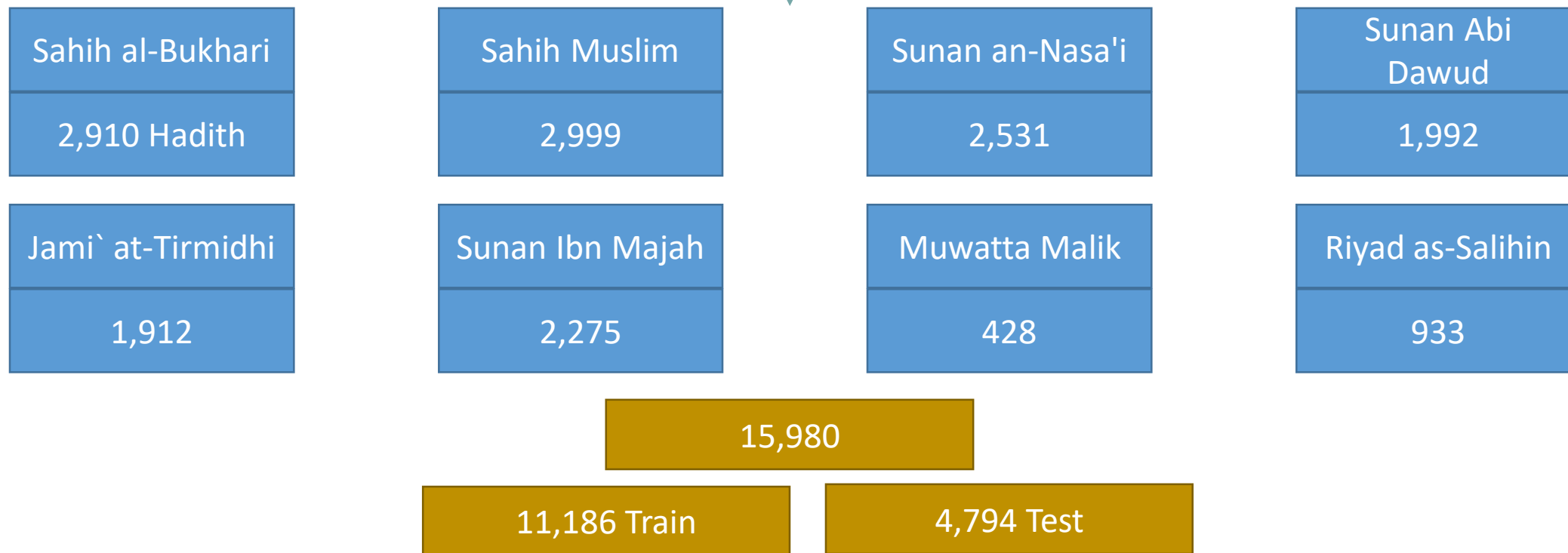
Target: Develop NLP Model to segment Hadith



Training Data



Scraped



Data Processing

حَدَّثَنِي عَبْدُ الْحَمِيدِ بْنُ بَيَانَ الْوَاسِطِيُّ، حَدَّثَنَا خَالِدٌ، - يَعْنِي ابْنَ عَبْدِ اللَّهِ - عَنْ سُهَيْلٍ، عَنْ أَبِيهِ، عَنْ أَبِي هُرَيْرَةَ، قَالَ قَالَ رَسُولُ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ " إِذَا أَدَّنَ الْمُؤَذِّنُ أَدْبَرَ الشَّيْطَانُ وَلَهُ حُصَاصٌ "

Remove Rubbish

- Eng. Alphanumeric
- Special Characters
- Special Unicode

Remove White Spaces

Remove Punctuation

Remove Tashkeel

حَدَّثَنِي عَبْدُ الْحَمِيدِ بْنُ بَيَانَ الْوَاسِطِيُّ حَدَّثَنَا خَالِدٌ يَعْنِي ابْنَ عَبْدِ اللَّهِ عَنْ سُهَيْلٍ عَنْ أَبِيهِ عَنْ أَبِي هُرَيْرَةَ قَالَ قَالَ رَسُولُ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ إِذَا أَدَّنَ الْمُؤَذِّنُ أَدْبَرَ الشَّيْطَانُ وَلَهُ حُصَاصٌ

Train Dataset

حَدَّثَنِي عَبْدُ الْحَمِيدِ بْنُ بَيَانَ الْوَاسِطِيُّ حَدَّثَنَا خَالِدٌ يَعْنِي ابْنَ عَبْدِ اللَّهِ عَنْ سُهَيْلٍ عَنْ أَبِيهِ عَنْ أَبِي هُرَيْرَةَ قَالَ قَالَ رَسُولُ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ

إِذَا أَدَّنَ الْمُؤَذِّنُ أَدْبَرَ الشَّيْطَانُ وَلَهُ حُصَاصٌ

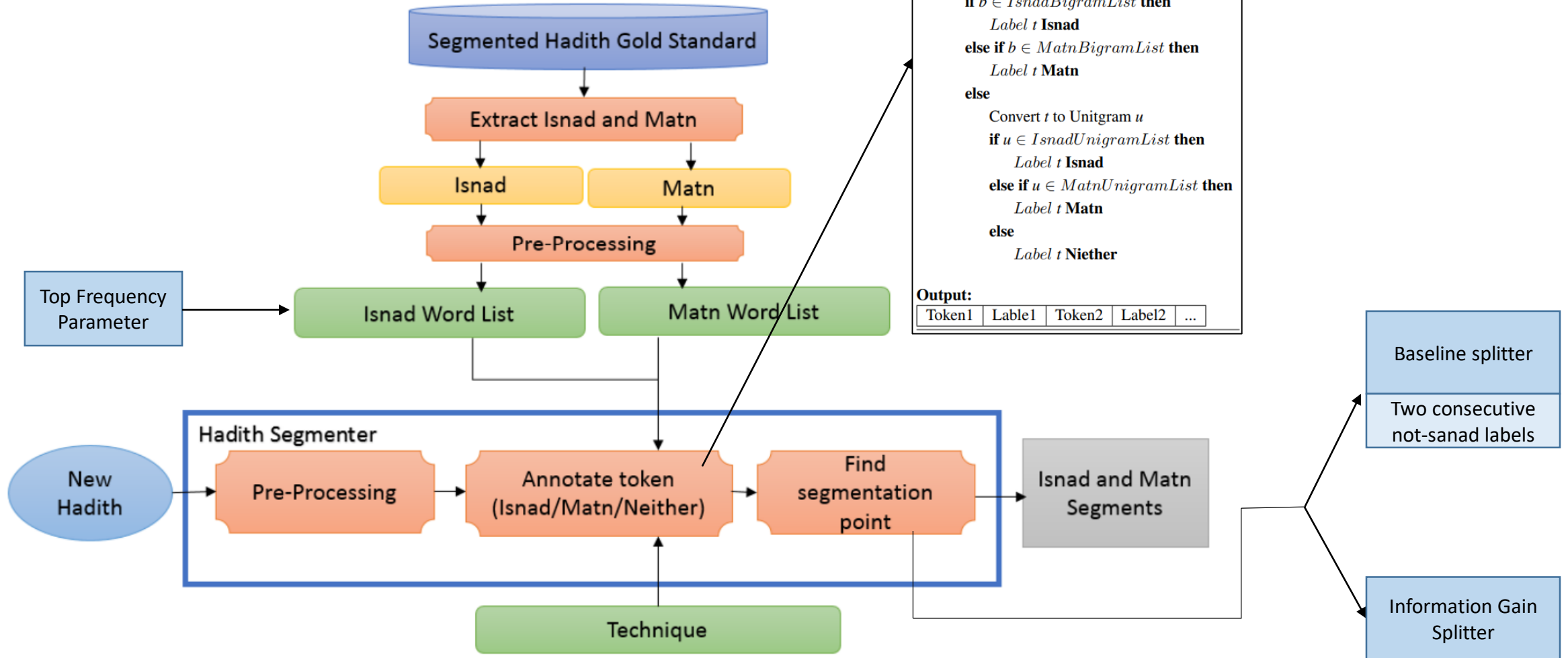
Test Dataset

حَدَّثَنِي عَبْدُ الْحَمِيدِ بْنُ بَيَانَ الْوَاسِطِيُّ حَدَّثَنَا خَالِدٌ يَعْنِي ابْنَ عَبْدِ اللَّهِ عَنْ سُهَيْلٍ عَنْ أَبِيهِ عَنْ أَبِي هُرَيْرَةَ قَالَ قَالَ رَسُولُ اللَّهِ صَلَّى اللَّهُ عَلَيْهِ وَسَلَّمَ إِذَا أَدَّنَ الْمُؤَذِّنُ أَدْبَرَ الشَّيْطَانُ وَلَهُ حُصَاصٌ

وسلم

27

Hadith Segmenter



Algorithm 1 Annotate Tri-gram tokens

Tokenize Hadith into Tri-grams T

for $t \in T$ **do**

if $t \in \text{IsnadTrigramList}$ **then**

 Label t **Isnad**

else if $t \in \text{MatnTrigramList}$ **then**

 Label t **Matn**

else

 Convert t to Bigrams b

if $b \in \text{IsnadBigramList}$ **then**

 Label t **Isnad**

else if $b \in \text{MatnBigramList}$ **then**

 Label t **Matn**

else

 Convert t to Unitgram u

if $u \in \text{IsnadUnigramList}$ **then**

 Label t **Isnad**

else if $u \in \text{MatnUnigramList}$ **then**

 Label t **Matn**

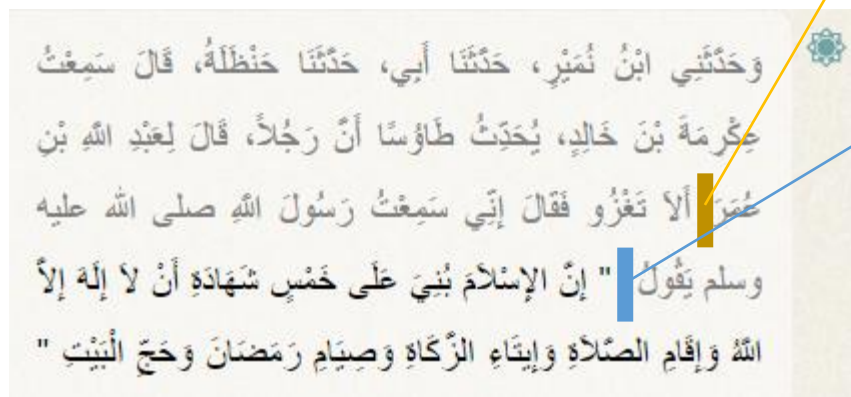
else

 Label t **Niether**

Output:

Token1	Label1	Token2	Label2	...
--------	--------	--------	--------	-----

Evaluation



Predicted Split Position

\widehat{spos}

21

Split Position

$spos$

33

$|33-21| = 12$

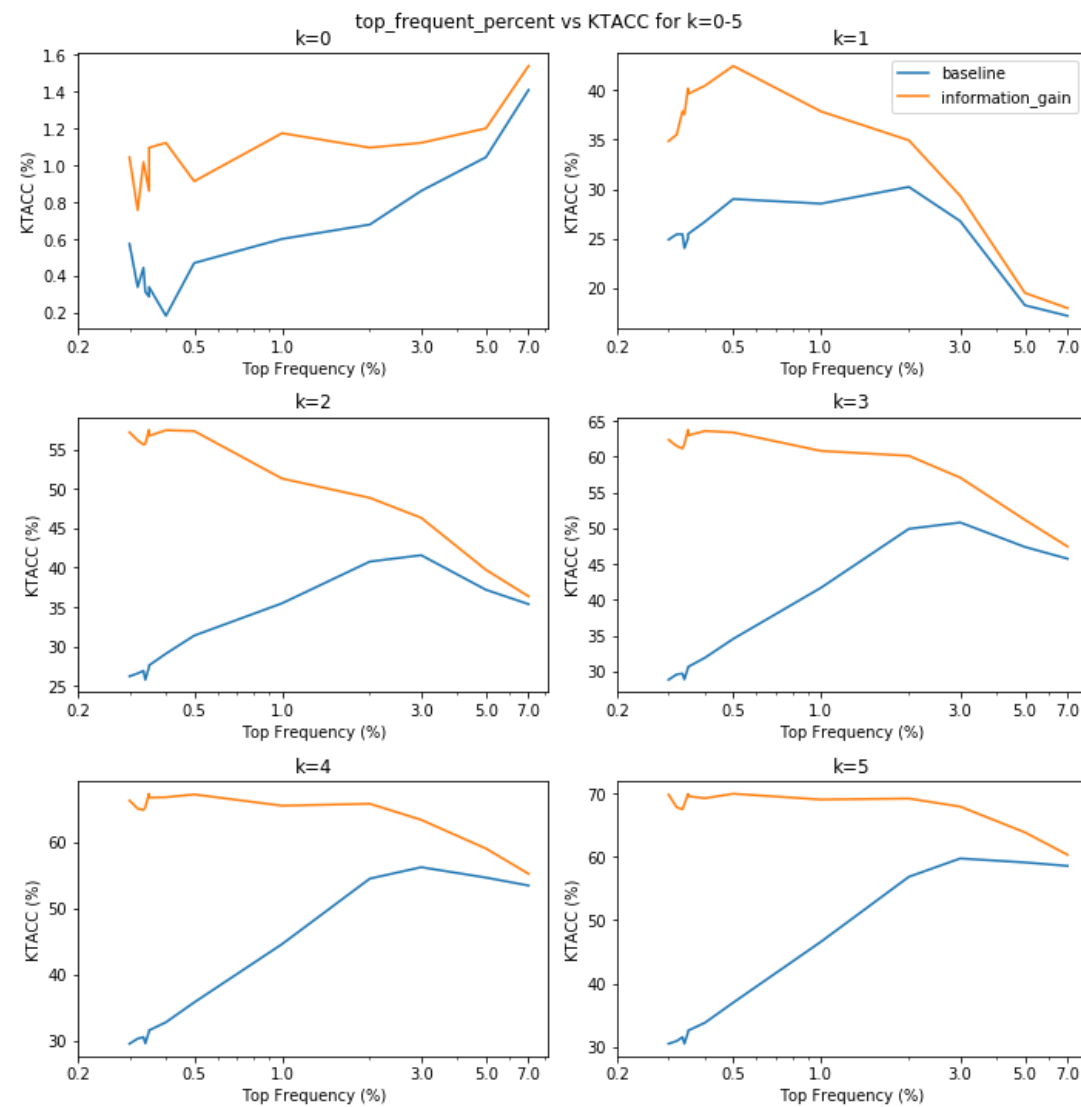
Compare with k

$$T_{pred}(hadith) = 1(\text{abs}(spos - \widehat{spos}) \leq k)$$

$$F_{pred}(hadith) = 1 - T_{pred}(hadith)$$

$$KTACC = \frac{\sum_{h=0}^H T_{pred}(h)}{\sum_{h=0}^H T_{pred}(h) + \sum_{h=0}^H F_{pred}(h)}$$

Fine Tuning



Thank you