

# Similarity Analysis via NLP Models

On Ahadis



# Process

1. Preprocessing of available data
2. Vectorization of data
3. Similarity Computation
4. Hadith Equivalence class Generation
5. Result Analysis



01

Preprocessing



# Understanding the Problem



Removing  
articles

Neptune is the  
farthest planet  
from the Sun



Removing Araab

أنا أحب أكل التفاح

...



Checking for  
Data  
Inconsistency

...



02

# Vectorization





# NLP Models



01

TDFIDF Model

02

BERT Model



# TFIDF Model

A word frequency based model

TF-IDF is an abbreviation for Term Frequency Inverse Document Frequency

Count Vectorizer

	blue	bright	sky	sun
Doc1	1	0	1	0
Doc2	0	1	0	1

TD-IDF Vectorizer

	blue	bright	sky	sun
Doc1	0.707107	0.000000	0.707107	0.000000
Doc2	0.000000	0.707107	0.000000	0.707107

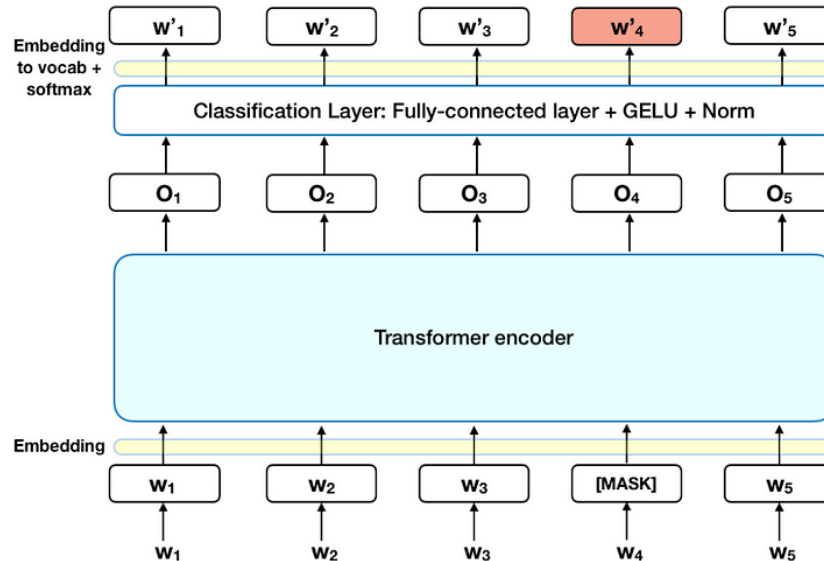


# BERT Model

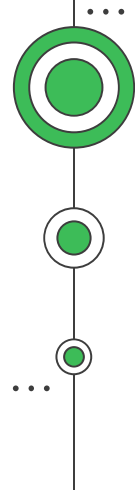
A word semantic based model

BERT makes use of Transformer, an attention mechanism that learns contextual relations between words (or sub-words) in a text.

ARA Vec for Arabic part







# 03

## Similarity Computation



# Chapter Wise

... كتاب صلاة التراويح

66 . كتاب الجمعة

67 . كتاب الجزية والموادعة

68 . كتاب الإجارة

69 . كتاب اللباس

70 . كتاب أخبار الآحاد

71 . كتاب أحاديث الأنبياء

72 . كتاب التوحيد

73 . كتاب الوتر

74 . كتاب المزارعة

75 . كتاب الشفعة

76 . كتاب المعطالم

77 . كتاب المرضى

78 . كتاب التقصير

79 . كتاب الوضوء

80 . كتاب المحصر

81 . كتاب الكفالي

82 . كتاب المساقاة

83 . كتاب النكاح

85 . كتاب العيدين

86 . كتاب الاعتكاف

87 . كتاب الحوالات

88 . كتاب فضائل المدينة

89 . كتاب الأذان

90 . كتاب مناقب الأنصار

91 . كتاب الإيمان

92 . كتاب الإكراه

93 . كتاب الاستئذان

94 . كتاب فرض الخمس

95 . كتاب الزكاة

96 . كتاب الجنائز

97 . كتاب الهبة وفضلها والتحريض عليها

98 . كتاب الأضاحي

99 . كتاب الطلاق

100 . كتاب الدعوات

101 . كتاب فضل الصلاة في مسجد مكة والمدينة

102 . كتاب السلم

103 . كتاب الشركة

## Distinct Strings:

1 . كتاب التمنى

2 . كتاب فضائل القرآن

3 . كتاب الرهن

4 . كتاب الكسوف

5 . كتاب الغسل

6 . كتاب كفارات الأيمان

7 . كتاب الوصايا

8 . كتاب الطب

9 . كتاب التعبير

استنابة المرتدين والمعاندين وقتالهم

11 . كتاب الشروط

12 . كتاب الحيل

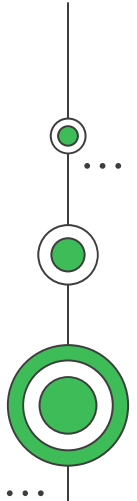
13 . كتاب القدر

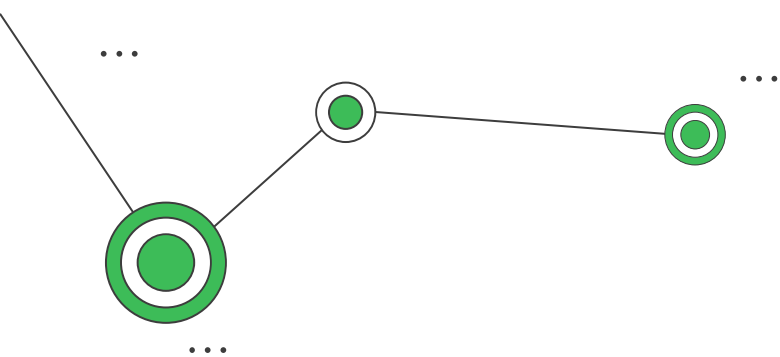
14 . كتاب سجود القرآن

15 . كتاب العتق

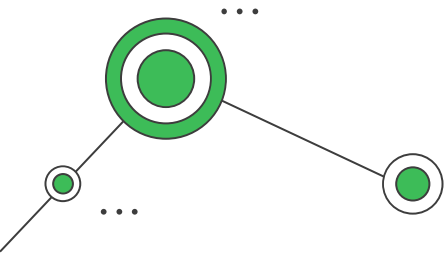
16 . كتاب الأيمان والندور

17 . كتاب الرقاق





# Whole Data





# Mukarrarat

Similar Ahadis Data taken from GitHub



04

# Equivalence Class



# Equivalence class Analysis

Identical : 0.9888-1  
Similar : 0.8- 0.988  
Not Similar: < 0.8

