# **Nugget**: Neural Agglomerative Embeddings of Text

Guanghui Qin      Benjamin Van Durme

JOHNS HOPKINS UNIVERSITY

**Solution**

**Extrinsic**

## Cram sentence into a vector? No!

Single vector   *limited capacity*

Token-wise   *non-scalable*

*I think , therefore I am .*

Nugget   *dynamic balance*

Encode text with ***variable*** number of vectors!

### What are used as nuggets?

Run Nugget scorer and label nugget with boxes:

NLP is an interdisciplinary subfield of linguistics , computer science , and artificial intelligence concerned with the interactions between computers and human language , in particular how to program computers to process and …

Surprise? Nuggets are *text delimiters*!

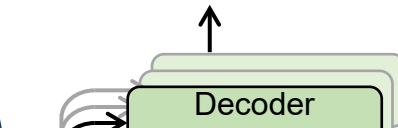And each nugget encodes *its preceding texts*!

## How?

*Key idea: Use the vectors of a **subset** of tokens to represent the whole passage*
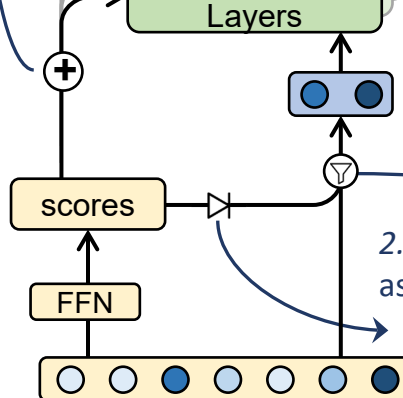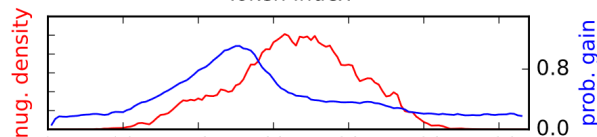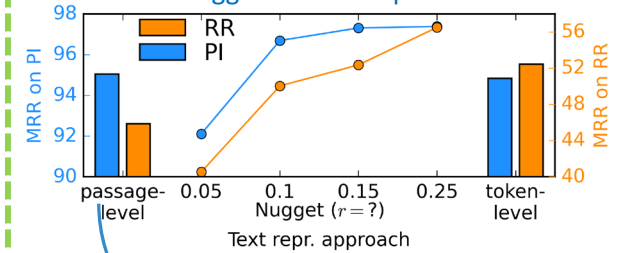
5. A residual connection from scores to the cross-attention to ensure differentiability

I think, therefore I am.

4. Train as *AutoEncoder* or *Machine Translation*

Decoder Layers

3. $k$ vectors are used as passage representation, called ***nuggets***

scores

*2.* A FFN scorer <u>selects</u> $k$ tokens as nuggets. Note selection op is <u>non-differentiable</u>.

FFN

Encoder

1. Apply a text encoder

I think , therefore I am .

**Intrinsic**

## What are they good for?

### Document similarity test

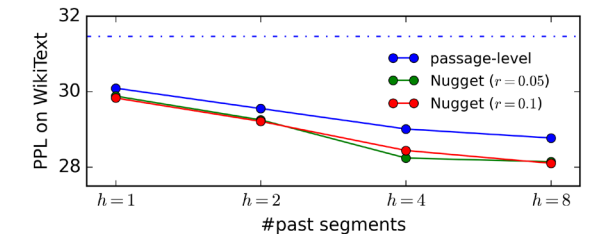Similarity between documents are measured by the many-to-many cosine similarity between their nugget vectors.

More nuggets -> better performance

RR   PI

MRR on PI: 98, 96, 94, 92, 90

MRR on RR: 56, 52, 48, 44, 40

passage-level   0.05   0.1   0.15   0.25   token-level

Nugget ($r = ?$)

Text repr. approach

Outperform the single-vector baseline

### LM with compressed context

Nugget is used to provide a compressed context to help a language model.

PPL on WikiText: 32, 30, 28

passage-level
Nugget ($r = 0.05$)
Nugget ($r = 0.1$)

$h = 1$   $h = 2$   $h = 4$   $h = 8$

#past segments

Nugget lowers the PPL. Longer compressed context consistently brings better performance.