

# Pre-trained Language Models Can be Fully Zero-Shot Learners

Xuandong Zhao , Siqi Ouyang , Zhiguo Yu , Ming Wu , Lei Li

ACL 2023

<https://aclanthology.org/2023.acl-long.869.pdf>

ACL読み会2023

名古屋大学 武田・笹野研究室

M1 大鹿 雅史

# 概要

- PLMをZero-shotでNLUタスクなどの性能を向上させる手法(NPPrompt)の提案
  - ・テキスト分類、感情分析など
- 入力プロンプトの[MASK]の単語予測とラベル名の近傍の単語を利用
- 他のZero-shot手法を大幅に上回る性能を示す

- 選定理由
  - ・手法がシンプルでわかりやすい
  - ・LLMによってZero-shotでの研究は増えてきそう

# 背景

- PLMの台頭によりNLUタスクの大きな進歩が見られる
  - ・ 事前学習によってコーパスから知識を獲得する
  - ・ タスク特有のデータセットでfine-tuningすることでNLPタスクの性能が向上する
- Fine-tuningの問題点
  - ・ ラベル付きデータセットの構築に人手のコスト
  - ・ モデルを学習させる計算コスト
  - ・ タスクごとにモデルを保存する必要

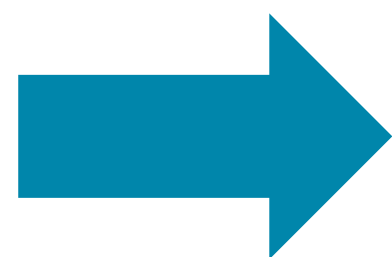
# 背景

## ■ PLMの台頭によりNLUタスクの大きな進歩が見られる

- ・ 事前学習によってコーパスから知識を獲得する
- ・ タスク特有のデータセットでfine-tuningすることでNLPタスクの性能が向上する

## ■ Fine-tuningの問題点

- ・ ラベル付きデータセットの構築に人手のコスト
- ・ モデルを学習させる計算コスト
- ・ タスクごとにモデルを保存する必要



Zero-shot でタスクが解けるとコストが少なく済む

# 関連研究

- Prompting
- Verbalizer Design
- Zero-shot Text Classification

# 関連研究

## ■ Prompting

- ・ GPT-3の成功によってより注目を集めている技術
- ・ さまざまな自然言語処理タスクで性能を示している

# 関連研究

## ■ Prompting

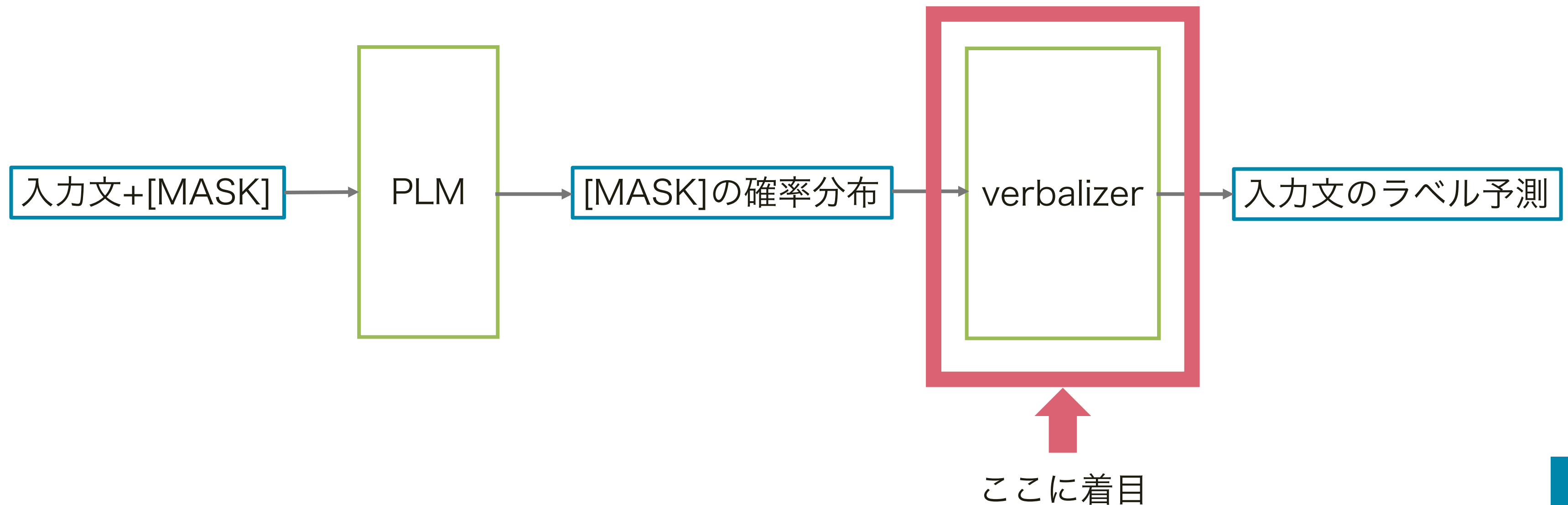
- ・ GPT-3の成功によってより注目を集めている技術
- ・ さまざまな自然言語処理タスクで性能を示している



# 関連研究

## ■ Prompting

- ・ GPT-3の成功によってより注目を集めている技術
- ・ さまざまな自然言語処理タスクで性能を示している





# 関連研究

## ■ Verbalizer Design

- ・ 人手で作成 → バイアスが発生
- ・ 自動でverbalizerを構築する手法が提案

### 1. [Auto-L](#)

ラベル名を表す単語の候補集合を元にモデルをfine-tuningしverbalizerを再定義する

### 2. [AutoPrompt](#)

勾配ベースの探索でラベルを表すverbalizer(トリガートークン)を定義

### 3. [KPT](#)

外部知識を利用してverbalizerを洗練する

# 関連研究

## ■ Verbalizer Design

- ・ 人手で作成 → バイアスが発生
- ・ 自動でverbalizerを構築する手法が提案

### 1. [Auto-L](#)

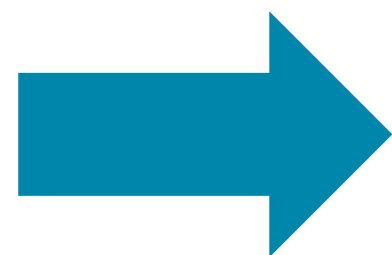
ラベル名を表す単語の候補集合を元にモデルをfine-tuningしverbalizerを再定義する

### 2. [AutoPrompt](#)

勾配ベースの探索でラベルを表すverbalizer(トリガートークン)を定義

### 3. [KPT](#)

外部知識を利用してverbalizerを洗練する



モデルの勾配情報や外部知識が必要

# 関連研究

## ■ Zero-shot Text Classification

- [SimPTC](#)

- 外部知識を元に各クラスのアンカーを決める
- アンカーとの類似度を元にクラスタリングを行いテキスト分類

- [LOTClass](#)

- 各クラスを表す単語を複数定義
- MLMの単語予測と各クラスの単語の一致度からテキスト分類
- クラス名をもとにMLMの自己学習を行う

# 関連研究

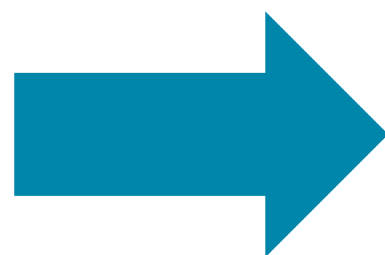
## ■ Zero-shot Text Classification

- [SimPTC](#)

- 外部知識を元に各クラスのアンカーを決める
- アンカーとの類似度を元にクラスタリングを行いテキスト分類

- [LOTClass](#)

- 各クラスを表す単語を複数定義
- MLMの単語予測と各クラスの単語の一致度からテキスト分類
- クラス名をもとにMLMの自己学習を行う

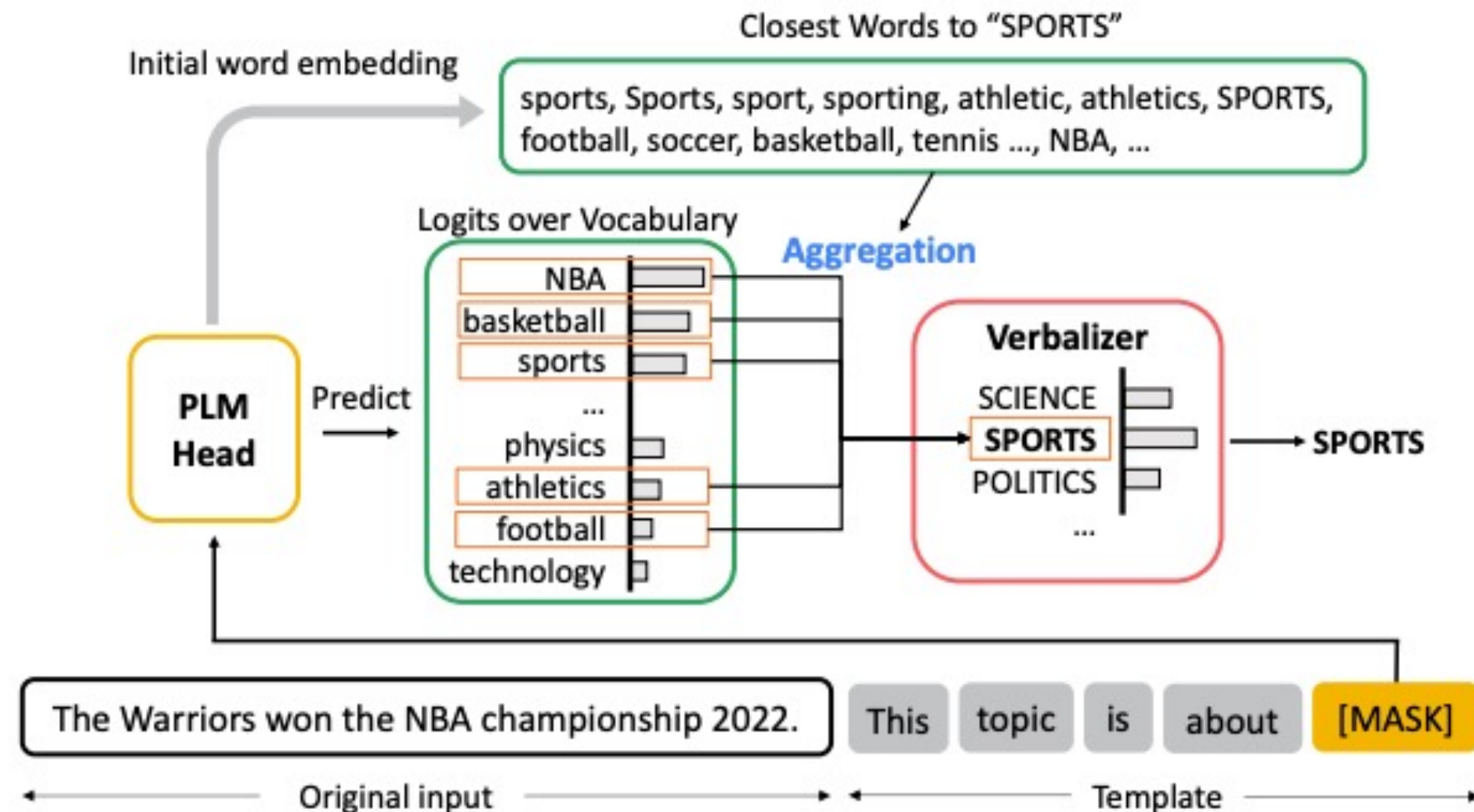


外部知識や学習コストがかかる

# NPPrompt : NonParametric Prompting

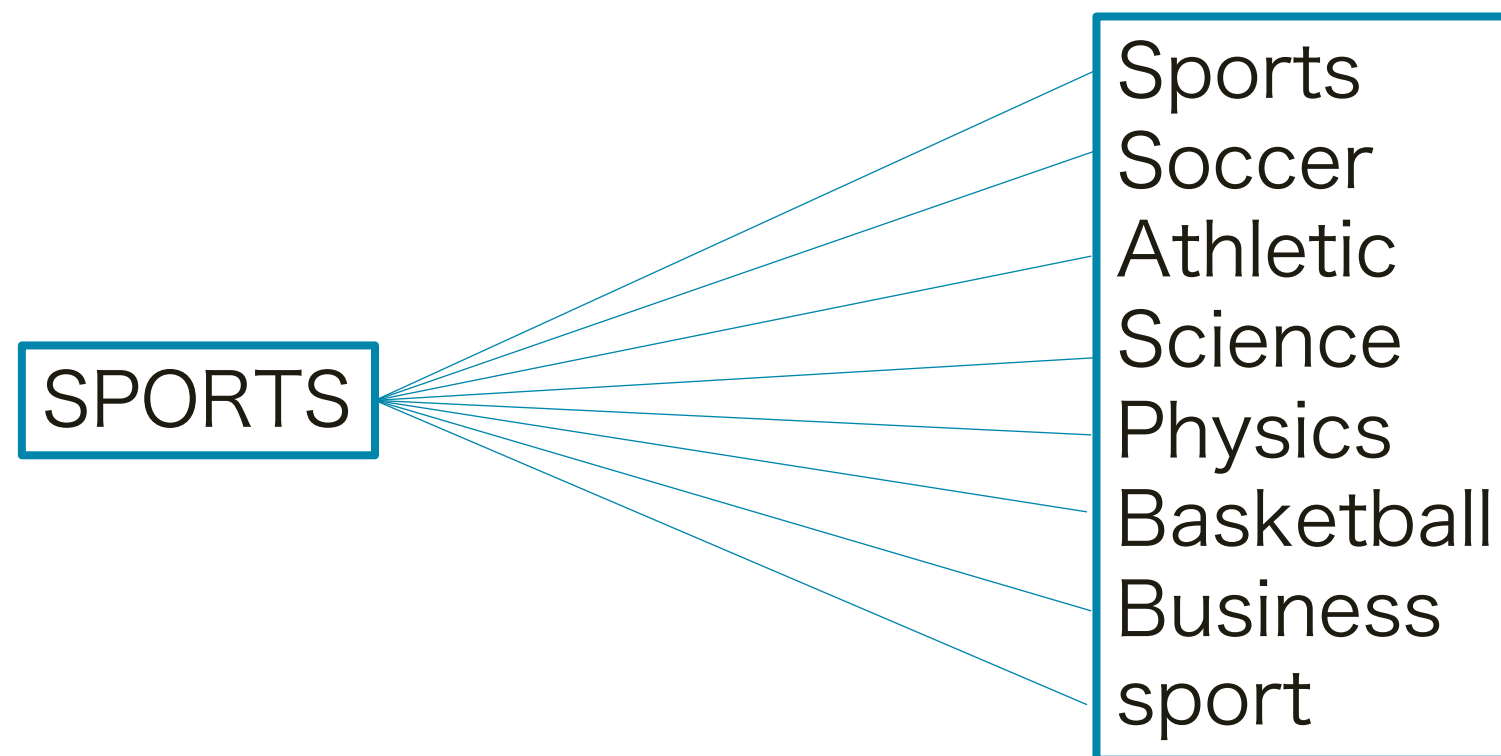
## ■ 手法

1. k-近傍verbalizerの構築
2. ラベルの予測



# NPPrompt : k-近傍verbalizerの構築

- PLMが持つ語彙とラベル名の埋め込みの類似度を計算
- 類似度が高いTop-kの単語からverbalizerを構築する

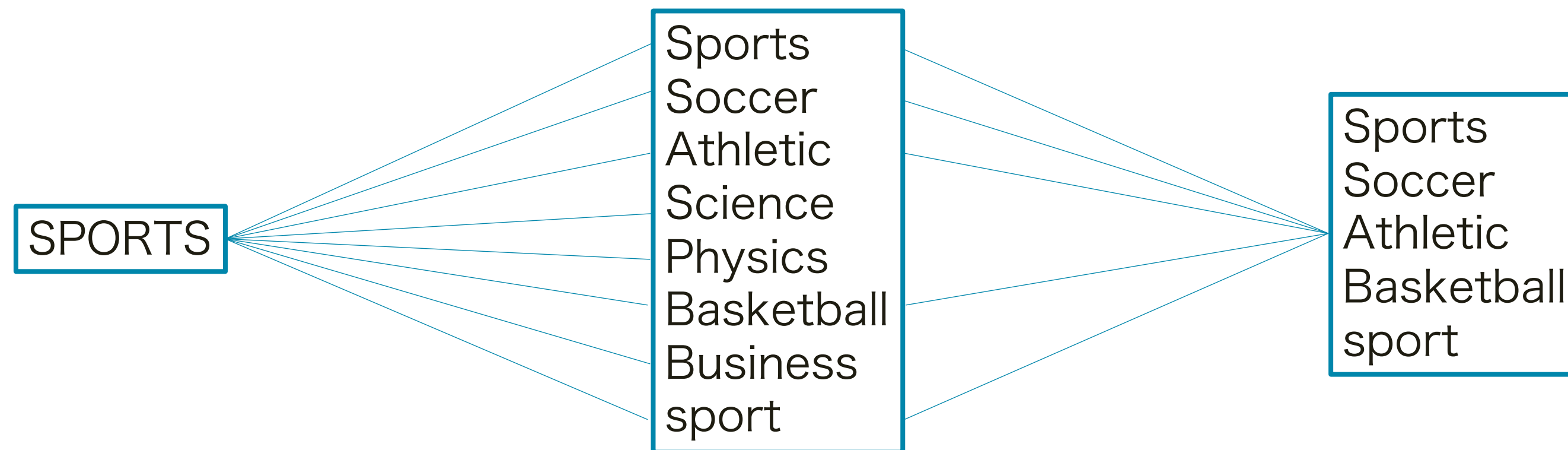


$$S(\text{emb}(v_i), \text{emb}(y_j)) = \frac{\text{emb}(v_i)}{\|\text{emb}(v_i)\|} \cdot \frac{\text{emb}(y_j)}{\|\text{emb}(y_j)\|}$$

$v_i$ : verbalizer内の単語,  $y_i$ : ラベル名

# NPPrompt : k-近傍verbalizerの構築

- PLMが持つ語彙とラベル名の埋め込みの類似度を計算
- 類似度が高いTop-kの単語からverbalizerを構築する



$$S(\text{emb}(v_i), \text{emb}(y_j)) = \frac{\text{emb}(v_i)}{\|\text{emb}(v_i)\|} \cdot \frac{\text{emb}(y_j)}{\|\text{emb}(y_j)\|} \quad M(y_j) = \text{Top } k\{S(\text{emb}(v_i), \text{emb}(y_j))\}_{v_i \in V}$$

$v_i$ : verbalizer内の単語,  $y_i$ : ラベル名

# NPPrompt : ラベルの予測

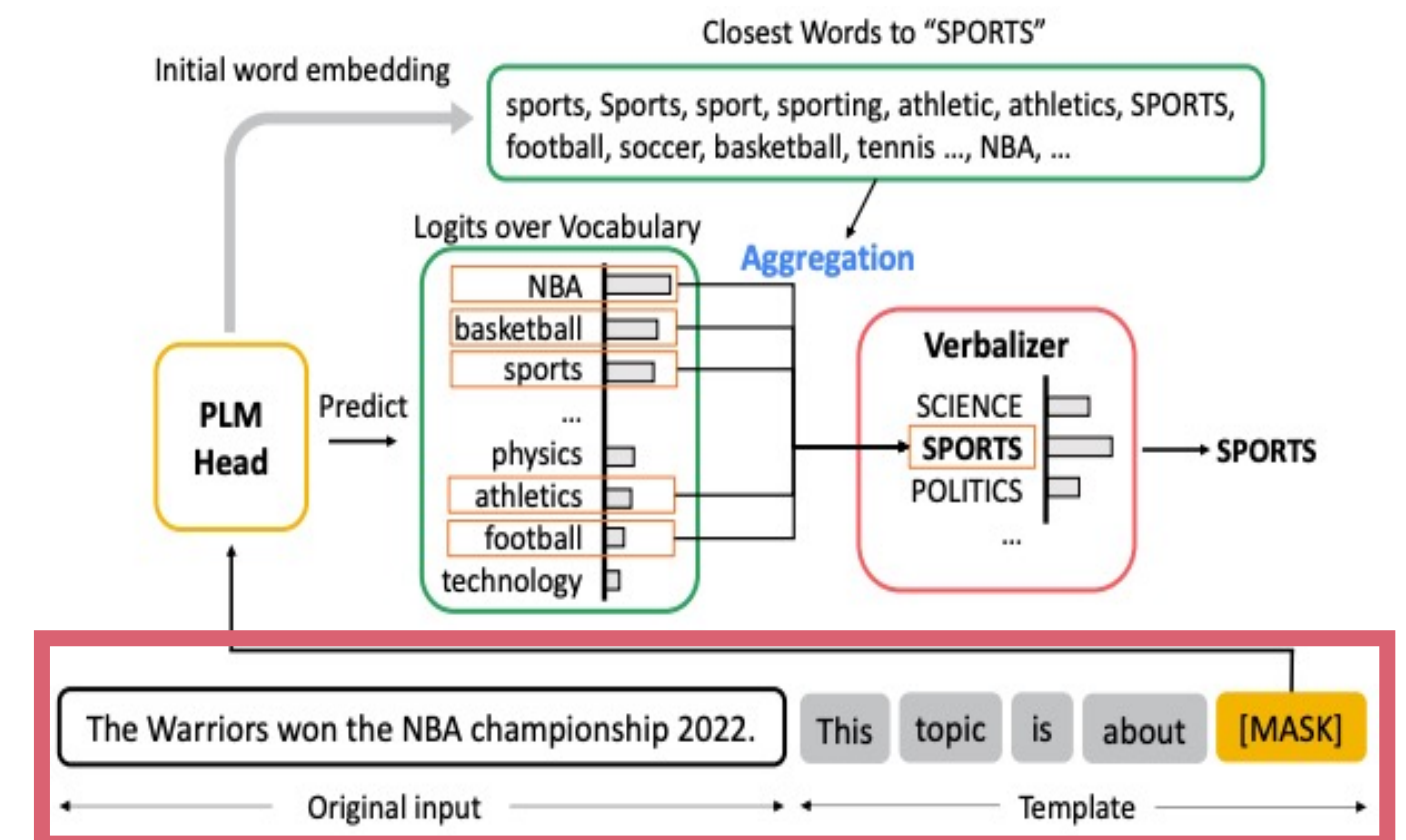
- 分類するテキストを元に[MASK]を利用したプロンプトを作成

The warriors won the NBA championship 2022.

→

- プロンプトをPLMに入力し[MASK]の予測Logitを出力

- Verbalizer内の単語からラベルを予測





# NPPrompt : ラベルの予測

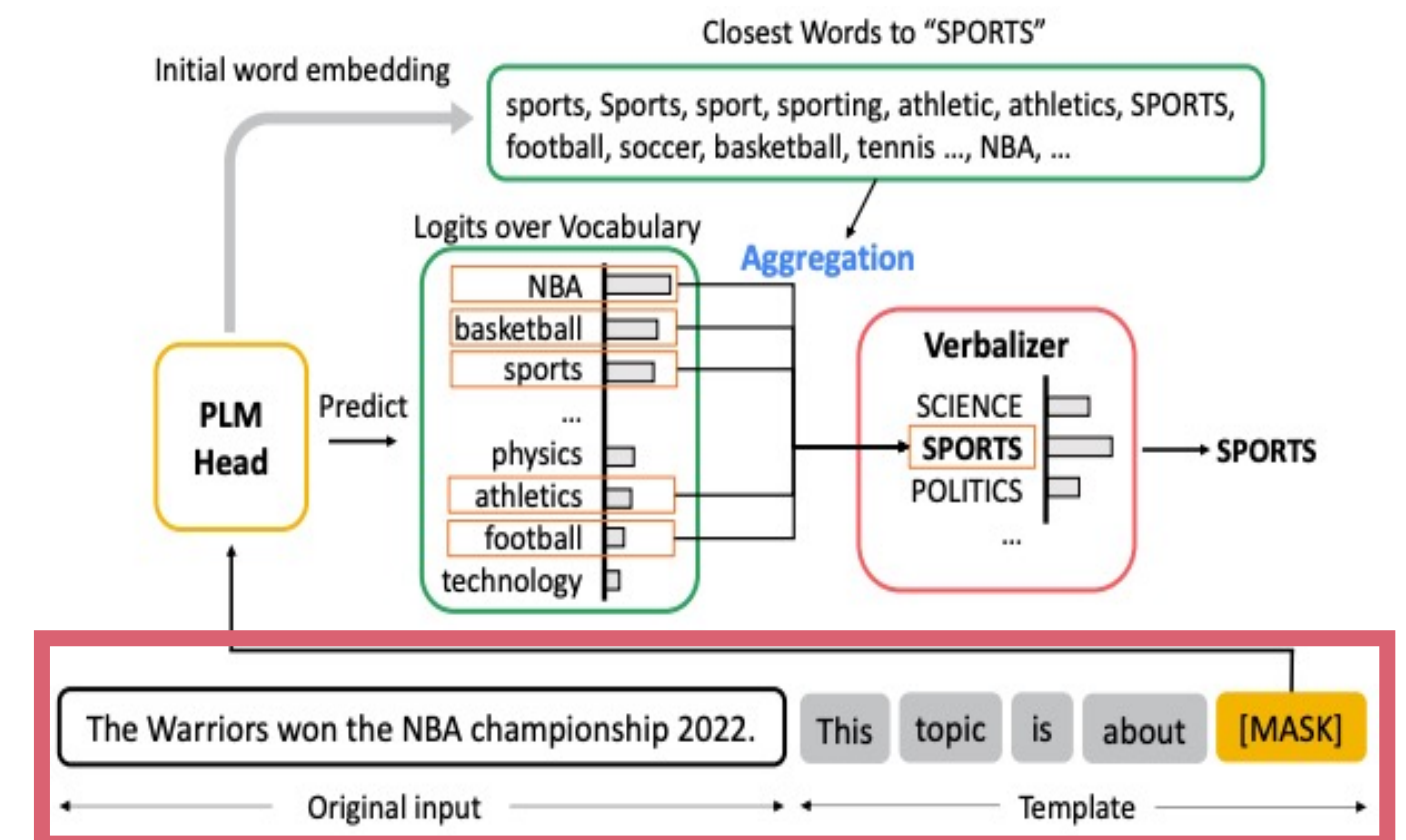
- 分類するテキストを元に[MASK]を利用したプロンプトを作成

The warriors won the NBA championship 2022.

→ The warriors won the NBA championship 2022. This topic is about [MASK].

- プロンプトをPLMに入力し[MASK]の予測Logitを出力

- Verbalizer内の単語からラベルを予測



# NPPrompt : ラベルの予測

- 分類するテキストを元に[MASK]を利用したプロンプトを作成

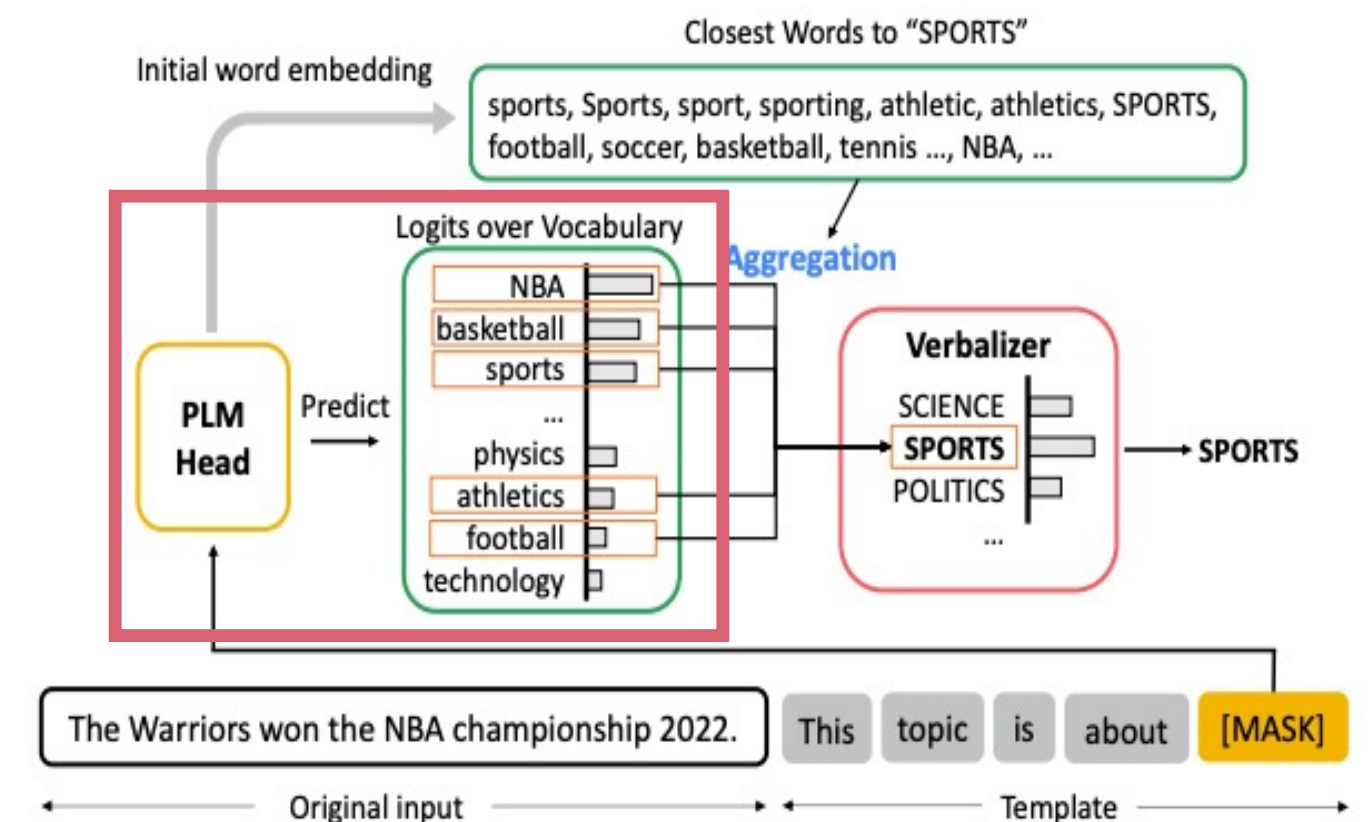
The warriors won the NBA championship 2022.

→ The warriors won the NBA championship 2022. This topic is about [MASK].

- プロンプトをPLMに入力し[MASK]の予測Logitを出力

$$\Theta([MASK] = v_i | x_{prompt} = \tau(x))$$

- Verbalizer内の単語からラベルを予測



# NPPrompt : ラベルの予測

- 分類するテキストを元に[MASK]を利用したプロンプトを作成

The warriors won the NBA championship 2022.

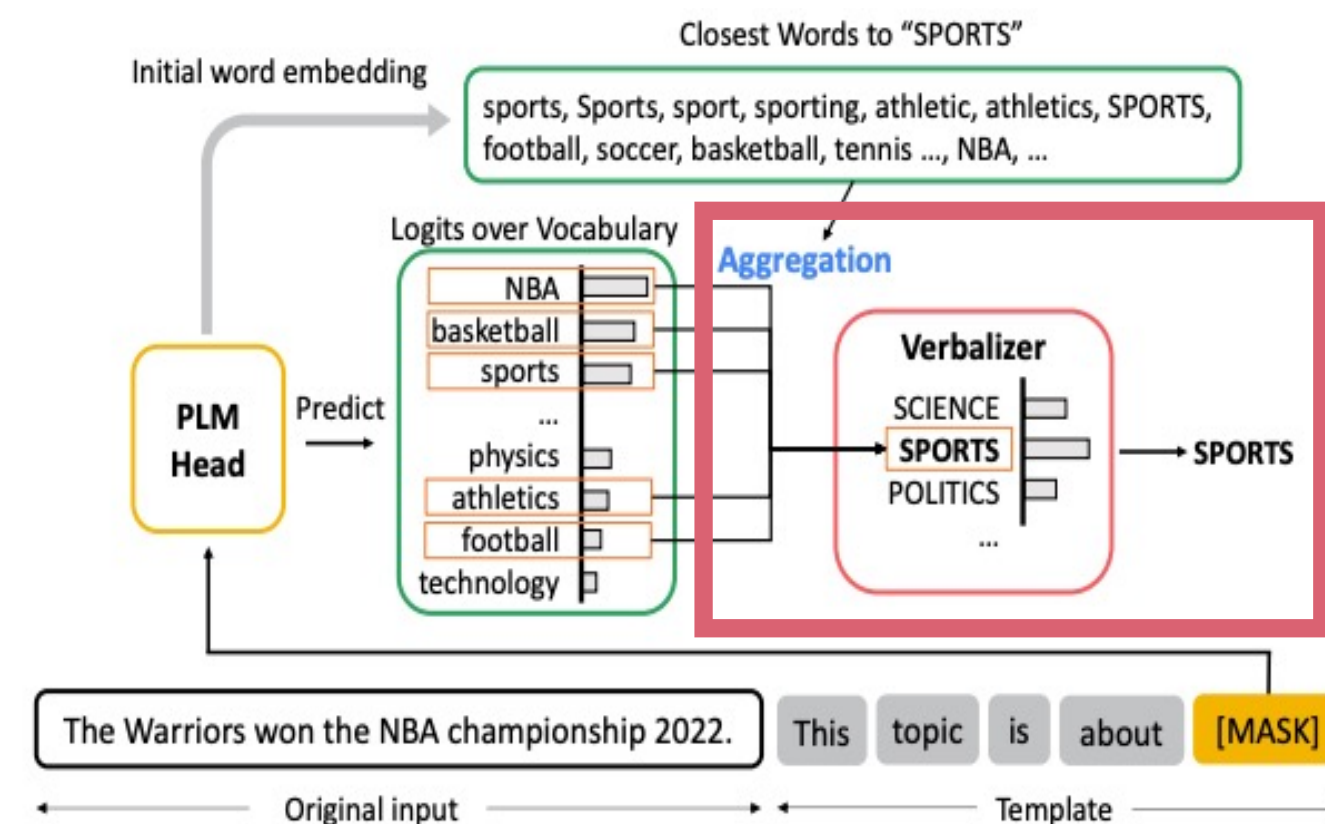
→ The warriors won the NBA championship 2022. This topic is about [MASK].

- プロンプトをPLMに入力し[MASK]の予測Logitを出力

$$\Theta([MASK] = v_i | x_{prompt} = \tau(x))$$

- Verbalizer内の単語からラベルを予測

$$Q(y_j | x) = \sum_{v_i \in M(y_j)} w(v_i, y_j) \cdot \Theta([MASK] = v_i | x_{prompt} = \tau(x))$$



# NPPrompt : ラベルの予測

- 分類するテキストを元に[MASK]を利用したプロンプトを作成

The warriors won the NBA championship 2022.

→ The warriors won the NBA championship 2022. This topic is about [MASK].

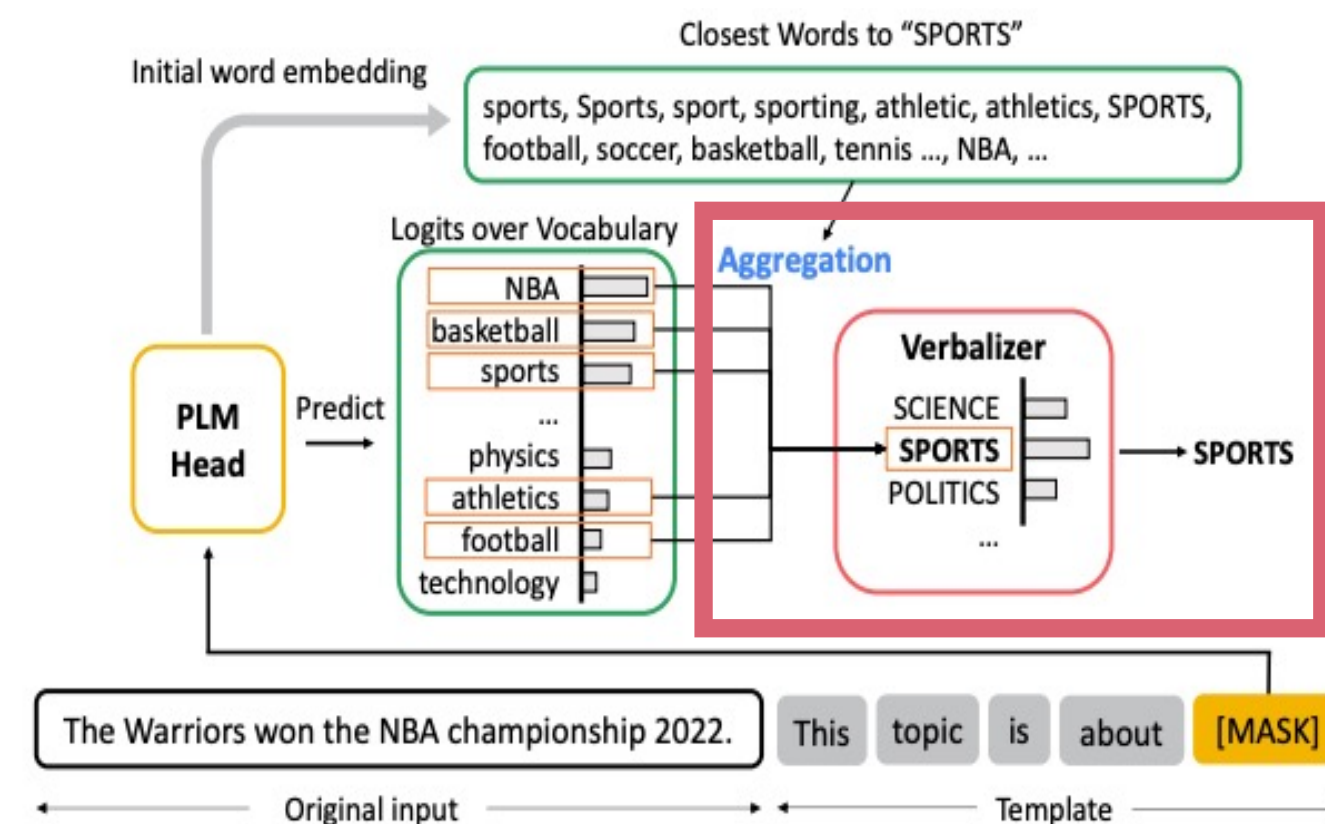
- プロンプトをPLMに入力し[MASK]の予測Logitを出力

$$\Theta([MASK] = v_i | x_{prompt} = \tau(x))$$

- Verbalizer内の単語からラベルを予測

$$Q(y_j | x) = \sum_{v_i \in M(y_j)} \underbrace{w(v_i, y_j)}_{\text{ソフトマックス関数}} \cdot \Theta([MASK] = v_i | x_{prompt} = \tau(x))$$

ソフトマックス関数



# 実験設定

## ■ 利用するデータセット

### ① テキスト分類

- ・ トピック分類

[AG News](#) : ニュースのトピック分類 (4クラス)

[DBpedia](#) : Wikipediaのトピック分類 (14クラス)

- ・ 感情分析

[IMDB](#) : 映画レビューの感情分析 (positive : negative)

[Amazon](#) : レビューの感情分析 (positive : negative)

### ② [GLUEベンチマーク](#) (自然言語推論、含意関係認識など)

## ■ 利用モデル : RoBERTa-large(355M)



# 実験設定

## ■ 利用したプロンプトやラベル名の例①

Dataset	Template	Label Names	$k$
AG News	A [MASK] news : $x$ .	category 1: <i>world, politics</i> category 2: <i>sports</i> category 3: <i>business</i> category 4: <i>technology, science</i>	12
DBPedia	$x_1$ $x_2$ In this sentence, $x_1$ is a [MASK] .	category 1: <i>company</i> category 2: <i>school</i> category 3: <i>artist</i> category 4: <i>sports</i> category 5: <i>politics, office</i> category 6: <i>transportation, car, bus, train</i> category 7: <i>building, construct, room, tower</i> category 8: <i>river, lake, mountain</i> category 9: <i>village</i> category 10: <i>animal, pet</i> category 11: <i>plant</i> category 12: <i>album</i> category 13: <i>film</i> category 14: <i>book, publication</i>	7
IMDB	$x$ All in all, it was [MASK] .	positive: <i>good</i> negative: <i>bad</i>	500
Amazon	$x$ All in all, it was [MASK] .	positive: <i>good</i> negative: <i>bad</i>	170

# 実験設定

## ■ 利用したプロンプトやラベル名の例②

SST-2	$x_1$ It was [MASK] .	positive: <i>great</i> negative: <i>terrible</i>	9
MNLI	$x_1$ ? [MASK] , $x_2$	entailment: <i>yes</i> neutral: <i>maybe</i> contradiction: <i>no</i>	4
MNLI-mm	$x_1$ ? [MASK] , $x_2$	entailment: <i>yes</i> neutral: <i>maybe</i> contradiction: <i>no</i>	4
QNLI	$x_1$ ? [MASK] , $x_2$	entailment: <i>Yes, Indeed, Overall</i> not_entailment: <i>No, Well, However</i>	3
RTE	$x_1$ ? [MASK] , $x_2$	entailment: <i>Yes</i> not_entailment: <i>No</i>	10
MRPC	$x_1$ [MASK] , $x_2$	equivalent: <i>Yes</i> not_equivalent: <i>No</i>	9
QQP	$x_1$ [MASK] , $x_2$	equivalent: <i>Yes</i> not_equivalent: <i>No</i>	9
CoLA	$x_1$ This is [MASK] .	grammatical: <i>true</i> not_grammatical: <i>wrong</i>	7

# 実験設定：比較手法①

## ■ ManualVerb

ドメイン知識を持った人間の専門家がverbalizerを人手で定義

## ■ Semantic Retrieval

入力文の埋め込みとラベル説明文の類似度を測る

## ■ NSP-BERT

テキスト分類を含意関係認識と置き換えNext Sentence Predictionを利用して予測

## ■ GPT-3/ChatGPT with descriptions

タスクの説明をプロンプトに加えモデルに入力

---

*AG News* :

[Descriptions] Definition: In this task, you are given a sentence. Your job is to classify the following sentence into one of the four different categories. The categories are: “politics”, “sports”, “business”, and “technology”. Input: [x]. Output:

---



## 実験設定：比較手法②

- SimPTC
- LOTClass
- KPT
- Null Prompt  
[MASK]を文末に付与し直接ラベル名を推定
- Multi-Null prompting  
[MASK]を文頭と文末に付与し直接ラベルを推定

完全にZero-shotの手法

# 実験結果①：テキスト分類

## ■ ゼロショットの手法において一番良い性能を記録

Human：人手の作業がある  
 KB：外部知識を利用  
 Unlabeled：ラベルのないコーパスを利用

Method	Human/KB	Unlabeled	AG News	DBPedia	IMDB	Amazon	Avg.
Manual Verb	✓	✗	79.6 <sub>0.6</sub>	71.7 <sub>1.1</sub>	92.0 <sub>0.7</sub>	87.3 <sub>0.4</sub>	82.7
Semantic Retrieval	✓	✗	73.1 <sub>1.2</sub>	78.6 <sub>0.8</sub>	64.8 <sub>1.3</sub>	59.4 <sub>0.7</sub>	69.0
NSP-BERT	✓	✗	77.4 <sub>0.6</sub>	64.7 <sub>5.3</sub>	72.8 <sub>1.1</sub>	72.7 <sub>3.9</sub>	71.9
GPT-3 w. descriptions	✓	✗	83.4	82.5	88.8	89.4	86.0
ChatGPT w. descriptions	✓	✗	83.8	92.0	92.7	<b>95.8</b>	91.1
SimPTC	✓	✗	<b>86.9</b> <sub>0.3</sub>	<b>93.2</b> <sub>1.0</sub>	91.0 <sub>0.0</sub>	93.9 <sub>0.0</sub>	<b>91.3</b>
LOTClass w/o. self train	✗	✓	82.2	86.0	80.2	85.3	83.4
LOTClass	✗	✓	86.4	91.1	86.5	91.6	88.9
KPT	✓	✓	86.7	87.4	<b>94.0</b>	94.6	90.7
Null Prompt	✗	✗	67.9 <sub>2.0</sub>	56.8 <sub>3.9</sub>	82.5 <sub>1.5</sub>	89.4 <sub>1.0</sub>	74.2
Multi-Null Prompt	✗	✗	68.2 <sub>1.8</sub>	67.6 <sub>1.8</sub>	86.6 <sub>0.6</sub>	86.2 <sub>2.7</sub>	77.2
NPPrompt	✗	✗	<b>85.2</b> <sub>0.5</sub>	<b>86.8</b> <sub>0.1</sub>	<b>94.2</b> <sub>0.2</sub>	<b>93.9</b> <sub>0.0</sub>	<b>90.0</b>

# 実験結果①：テキスト分類

- 多くのベースラインより上回る結果が得られたが、少しだけ最高性能を下回った

Method	Human/KB	Unlabeled	AG News	DBPedia	IMDB	Amazon	Avg.
Manual Verb	✓	✗	79.6 <sub>0.6</sub>	71.7 <sub>1.1</sub>	92.0 <sub>0.7</sub>	87.3 <sub>0.4</sub>	82.7
Semantic Retrieval	✓	✗	73.1 <sub>1.2</sub>	78.6 <sub>0.8</sub>	64.8 <sub>1.3</sub>	59.4 <sub>0.7</sub>	69.0
NSP-BERT	✓	✗	77.4 <sub>0.6</sub>	64.7 <sub>5.3</sub>	72.8 <sub>1.1</sub>	72.7 <sub>3.9</sub>	71.9
GPT-3 w. descriptions	✓	✗	83.4	82.5	88.8	89.4	86.0
ChatGPT w. descriptions	✓	✗	83.8	92.0	92.7	<b>95.8</b>	91.1
SimPTC	✓	✗	<b>86.9</b> <sub>0.3</sub>	<b>93.2</b> <sub>1.0</sub>	91.0 <sub>0.0</sub>	93.9 <sub>0.0</sub>	<b>91.3</b>
LOTClass w/o. self train	✗	✓	82.2	86.0	80.2	85.3	83.4
LOTClass	✗	✓	86.4	91.1	86.5	91.6	88.9
KPT	✓	✓	86.7	87.4	<b>94.0</b>	94.6	90.7
Null Prompt	✗	✗	67.9 <sub>2.0</sub>	56.8 <sub>3.9</sub>	82.5 <sub>1.5</sub>	89.4 <sub>1.0</sub>	74.2
Multi-Null Prompt	✗	✗	68.2 <sub>1.8</sub>	67.6 <sub>1.8</sub>	86.6 <sub>0.6</sub>	86.2 <sub>2.7</sub>	77.2
NPPrompt	✗	✗	<b>85.2</b> <sub>0.5</sub>	<b>86.8</b> <sub>0.1</sub>	<b>94.2</b> <sub>0.2</sub>	<b>93.9</b> <sub>0.0</sub>	<b>90.0</b>



## 実験結果①：テキスト分類

- GPT-3(175B)のような大規模なモデルよりも小さいモデルで高い性能を示した

Method	Human/KB	Unlabeled	AG News	DBPedia	IMDB	Amazon	Avg.
Manual Verb	✓	✗	79.6 <sub>0.6</sub>	71.7 <sub>1.1</sub>	92.0 <sub>0.7</sub>	87.3 <sub>0.4</sub>	82.7
Semantic Retrieval	✓	✗	73.1 <sub>1.2</sub>	78.6 <sub>0.8</sub>	64.8 <sub>1.3</sub>	59.4 <sub>0.7</sub>	69.0
NSP-BERT	✓	✗	77.4 <sub>0.6</sub>	64.7 <sub>5.3</sub>	72.8 <sub>1.1</sub>	72.7 <sub>3.9</sub>	71.9
GPT-3 w. descriptions	✓	✗	83.4	82.5	88.8	89.4	86.0
ChatGPT w. descriptions	✓	✗	83.8	92.0	92.7	<b>95.8</b>	91.1
SimPTC	✓	✗	<b>86.9</b> <sub>0.3</sub>	<b>93.2</b> <sub>1.0</sub>	91.0 <sub>0.0</sub>	93.9 <sub>0.0</sub>	<b>91.3</b>
LOTClass w/o. self train	✗	✓	82.2	86.0	80.2	85.3	83.4
LOTClass	✗	✓	86.4	91.1	86.5	91.6	88.9
KPT	✓	✓	86.7	87.4	<b>94.0</b>	94.6	90.7
Null Prompt	✗	✗	67.9 <sub>2.0</sub>	56.8 <sub>3.9</sub>	82.5 <sub>1.5</sub>	89.4 <sub>1.0</sub>	74.2
Multi-Null Prompt	✗	✗	68.2 <sub>1.8</sub>	67.6 <sub>1.8</sub>	86.6 <sub>0.6</sub>	86.2 <sub>2.7</sub>	77.2
NPPrompt	✗	✗	<b>85.2</b> <sub>0.5</sub>	<b>86.8</b> <sub>0.1</sub>	<b>94.2</b> <sub>0.2</sub>	<b>93.9</b> <sub>0.0</sub>	<b>90.0</b>

## 実験結果②：GLUEベンチマーク

- Zero-shotの手法の性能を上回る性能を示した

	<b>MNLI</b> (acc)	<b>MNLI-mm</b> (acc)	<b>SST-2</b> (acc)	<b>QNLI</b> (acc)	<b>RTE</b> (acc)	<b>MRPC</b> (F1)	<b>QQP</b> (F1)	<b>CoLA</b> (Matt.)	<b>Avg.</b>
<i>With human designed prompts / few-shot data</i>									
Manual Label	50.8	51.7	83.6	50.8	51.3	61.9	49.7	2.0	50.2
In-context learning	<b>52.0</b> <sub>0.7</sub>	<b>53.4</b> <sub>0.6</sub>	84.8 <sub>1.3</sub>	53.8 <sub>0.4</sub>	60.4 <sub>1.4</sub>	45.7 <sub>6.0</sub>	36.1 <sub>5.2</sub>	-1.5 <sub>2.4</sub>	48.1
Auto-L	41.6 <sub>5.4</sub>	42.3 <sub>6.2</sub>	84.3 <sub>3.3</sub>	57.9 <sub>3.9</sub>	<b>61.9</b> <sub>7.5</sub>	<b>67.7</b> <sub>7.9</sub>	55.5 <sub>5.0</sub>	1.2 <sub>4.8</sub>	51.6
AMuLaP	50.8 <sub>2.1</sub>	52.3 <sub>1.8</sub>	<b>86.9</b> <sub>1.6</sub>	53.1 <sub>2.8</sub>	58.9 <sub>7.9</sub>	56.3 <sub>5.0</sub>	60.2 <sub>2.7</sub>	2.3 <sub>1.4</sub>	52.6
Few-shot fine-tuning	45.8 <sub>6.4</sub>	47.8 <sub>6.8</sub>	81.4 <sub>3.8</sub>	<b>60.2</b> <sub>6.5</sub>	54.4 <sub>3.9</sub>	76.6 <sub>2.5</sub>	<b>60.7</b> <sub>4.3</sub>	<b>33.9</b> <sub>14.3</sub>	<b>57.6</b>
<i>Fully zero-shot</i>									
Majority	32.7	33.0	50.9	49.5	52.7	<b>81.2</b>	0.0	0.0	37.5
Null Prompt	33.1 <sub>0.4</sub>	33.8 <sub>0.5</sub>	79.1 <sub>4.0</sub>	50.7 <sub>0.1</sub>	47.2 <sub>0.6</sub>	12.9 <sub>7.0</sub>	1.3 <sub>1.0</sub>	-1.1 <sub>2.0</sub>	32.1
Multi-Null Prompt	38.0 <sub>3.5</sub>	38.5 <sub>4.1</sub>	70.2 <sub>7.7</sub>	52.2 <sub>1.7</sub>	53.0 <sub>2.2</sub>	19.9 <sub>8.7</sub>	25.5 <sub>13.4</sub>	<b>6.2</b> <sub>2.0</sub>	37.9
NPPrompt	<b>45.7</b> <sub>0.6</sub>	<b>45.9</b> <sub>0.5</sub>	<b>86.3</b> <sub>1.2</sub>	<b>57.6</b> <sub>0.7</sub>	<b>55.0</b> <sub>3.4</sub>	79.8 <sub>1.6</sub>	<b>52.4</b> <sub>0.4</sub>	4.9 <sub>4.1</sub>	<b>53.5</b>

## 実験結果②：GLUEベンチマーク

- タスクの平均で見るとfew-shotの手法も上回る性能

	MNLI (acc)	MNLI-mm (acc)	SST-2 (acc)	QNLI (acc)	RTE (acc)	MRPC (F1)	QQP (F1)	CoLA (Matt.)	Avg.
<i>With human designed prompts / few-shot data</i>									
Manual Label	50.8	51.7	83.6	50.8	51.3	61.9	49.7	2.0	50.2
In-context learning	<b>52.0</b> <sub>0.7</sub>	<b>53.4</b> <sub>0.6</sub>	84.8 <sub>1.3</sub>	53.8 <sub>0.4</sub>	60.4 <sub>1.4</sub>	45.7 <sub>6.0</sub>	36.1 <sub>5.2</sub>	-1.5 <sub>2.4</sub>	48.1
Auto-L	41.6 <sub>5.4</sub>	42.3 <sub>6.2</sub>	84.3 <sub>3.3</sub>	57.9 <sub>3.9</sub>	<b>61.9</b> <sub>7.5</sub>	<b>67.7</b> <sub>7.9</sub>	55.5 <sub>5.0</sub>	1.2 <sub>4.8</sub>	51.6
AMuLaP	50.8 <sub>2.1</sub>	52.3 <sub>1.8</sub>	<b>86.9</b> <sub>1.6</sub>	53.1 <sub>2.8</sub>	58.9 <sub>7.9</sub>	56.3 <sub>5.0</sub>	60.2 <sub>2.7</sub>	2.3 <sub>1.4</sub>	52.6
Few-shot fine-tuning	45.8 <sub>6.4</sub>	47.8 <sub>6.8</sub>	81.4 <sub>3.8</sub>	<b>60.2</b> <sub>6.5</sub>	54.4 <sub>3.9</sub>	76.6 <sub>2.5</sub>	<b>60.7</b> <sub>4.3</sub>	<b>33.9</b> <sub>14.3</sub>	<b>57.6</b>
<i>Fully zero-shot</i>									
Majority	32.7	33.0	50.9	49.5	52.7	<b>81.2</b>	0.0	0.0	37.5
Null Prompt	33.1 <sub>0.4</sub>	33.8 <sub>0.5</sub>	79.1 <sub>4.0</sub>	50.7 <sub>0.1</sub>	47.2 <sub>0.6</sub>	12.9 <sub>7.0</sub>	1.3 <sub>1.0</sub>	-1.1 <sub>2.0</sub>	32.1
Multi-Null Prompt	38.0 <sub>3.5</sub>	38.5 <sub>4.1</sub>	70.2 <sub>7.7</sub>	52.2 <sub>1.7</sub>	53.0 <sub>2.2</sub>	19.9 <sub>8.7</sub>	25.5 <sub>13.4</sub>	<b>6.2</b> <sub>2.0</sub>	37.9
NPPrompt	<b>45.7</b> <sub>0.6</sub>	<b>45.9</b> <sub>0.5</sub>	<b>86.3</b> <sub>1.2</sub>	<b>57.6</b> <sub>0.7</sub>	<b>55.0</b> <sub>3.4</sub>	79.8 <sub>1.6</sub>	<b>52.4</b> <sub>0.4</sub>	4.9 <sub>4.1</sub>	<b>53.5</b>

# 考察

1. 重み付け関数と類似度関数による性能差
2. ロジットと確率の性能差
3. k-近傍の単語数による性能差
4. 言語モデルによる性能差
5. テキスト分類以外へのタスクへの適用



# 考察①：重み付け関数と類似度関数による違い

- 類似度関数をコサイン類似度に固定し重み付け関数を変更

$$Q(y_j|x) = \sum_{v_i \in M(y_j)} w(v_i, y_j) \cdot \Theta([MASK] = v_i | x_{prompt} = \tau(x))$$

$$1. \textit{Default} \quad = w(v_i, y_j) = \frac{\exp(S(\textit{emp}(v_i), \textit{emb}(y_j)))}{\sum_{v_i \in M(y_j)} \exp(S(\textit{emp}(v_i), \textit{emb}(y_j)))}$$

$$2. \textit{Same weight} \quad = w(v_i, y_j) = 1$$

$$3. \textit{average weight} = w(v_i, y_j) = \frac{S(\textit{emp}(v_i), \textit{emb}(y_j))}{\sum_{v_i \in M(y_j)} S(\textit{emp}(v_i), \textit{emb}(y_j))}$$



# 考察①：重み付け関数と類似度関数による違い

■ 重み付け関数による大きな性能差のない結果

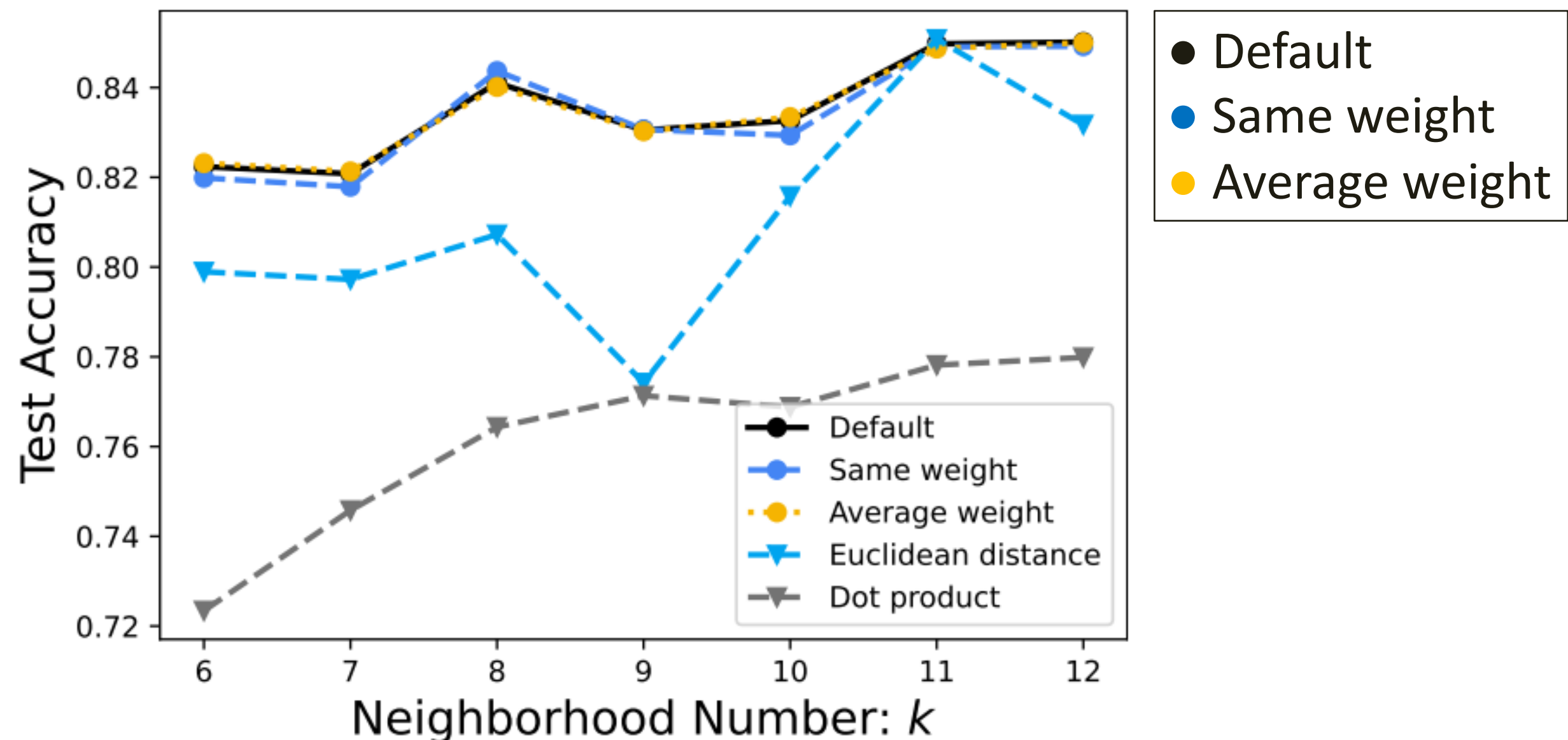


Figure 2: Effects of different aggregation.

## 考察①：重み付け関数と類似度関数による違い

- 重み付け関数を  $w(v_i, y_j) = 1$  に固定し類似度関数を変更

$$1. \textit{ cosine similarity} = S(\textit{emb}(v_i), \textit{emb}(y_i)) = \frac{\textit{emb}(v_i)}{\|\textit{emb}(v_i)\|} \cdot \frac{\textit{emb}(y_i)}{\|\textit{emb}(y_i)\|}$$

$$2. \textit{ Euclidean distance} = S(\textit{emb}(v_i), \textit{emb}(y_i)) = -\|\textit{emb}(v_i) - \textit{emb}(y_j)\|$$

$$3. \textit{ Dot product} = S(\textit{emb}(v_i), \textit{emb}(y_i)) = \textit{emb}(v_i) \cdot \textit{emb}(y_j)$$

# 考察①：重み付け関数と類似度関数による違い

■ 類似度関数によっては大きく性能に差が出る結果

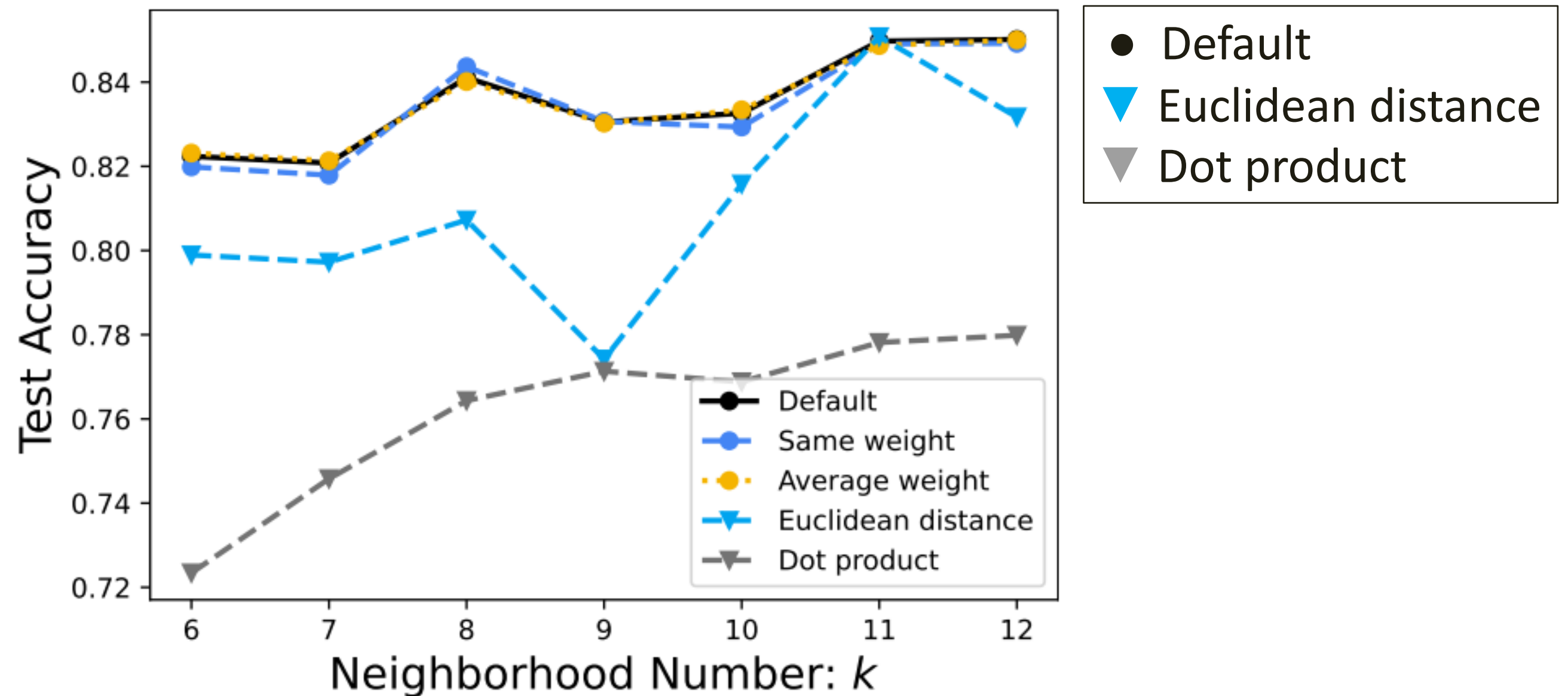


Figure 2: Effects of different aggregation.

## 考察②：ロジットと確率の性能差

- 近傍の単語数が少ない場合”sum logit”、多い場合は”sum prob”の性能が良い
- 一番高いスコアはsum logitの  $k=12$  の時

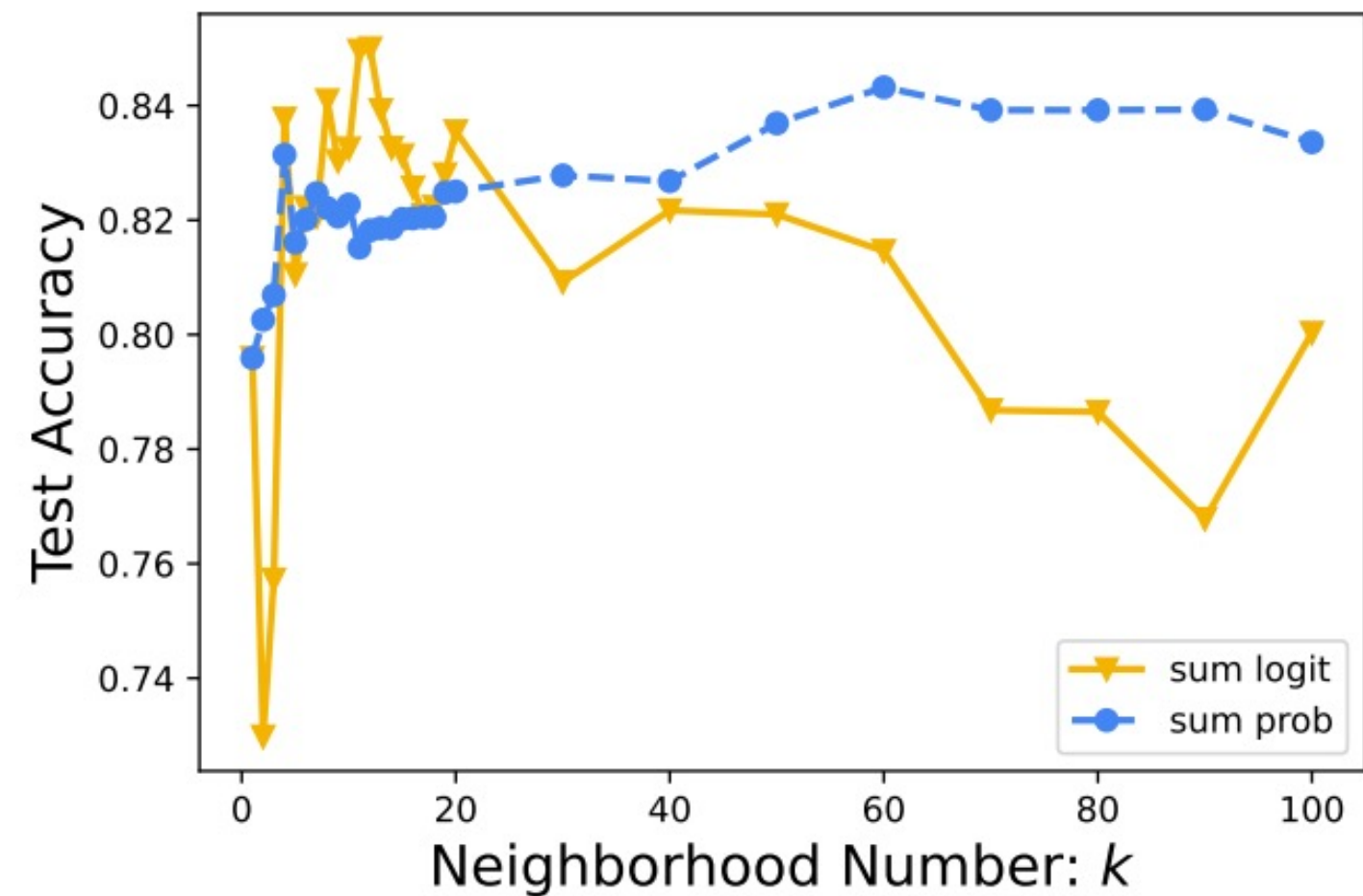
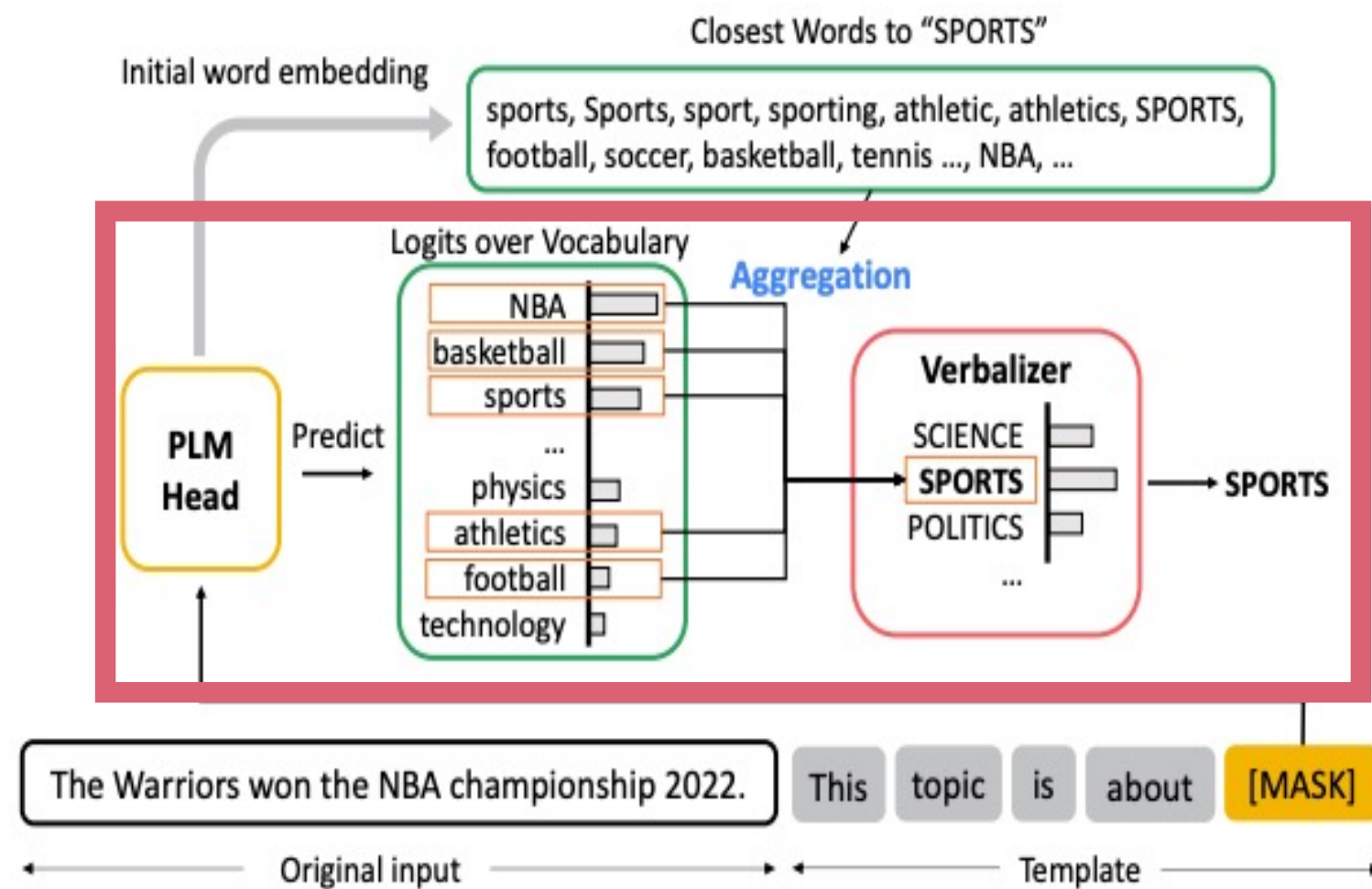
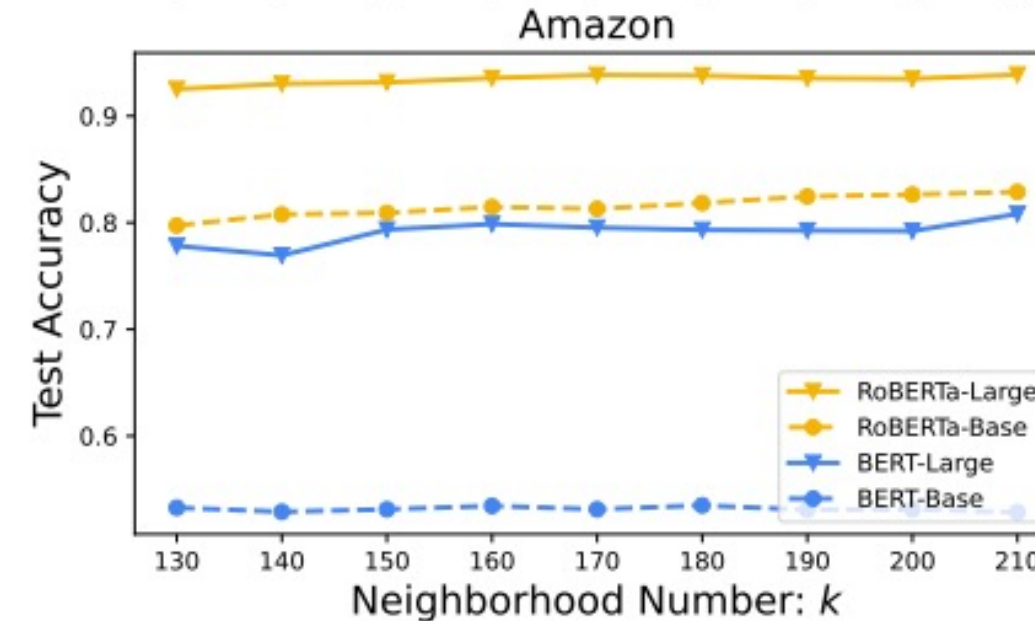
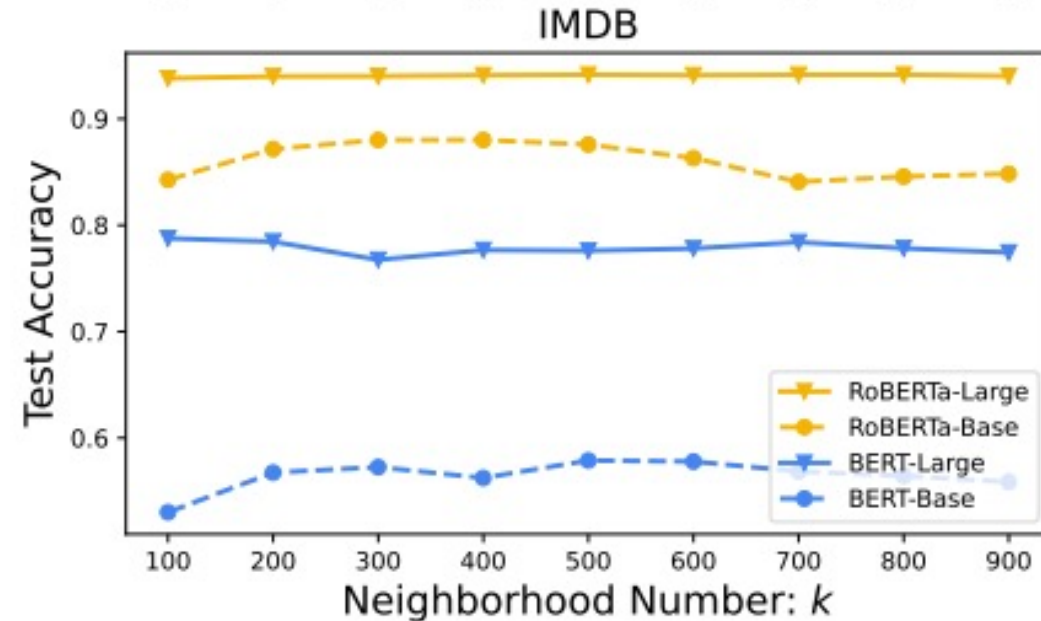
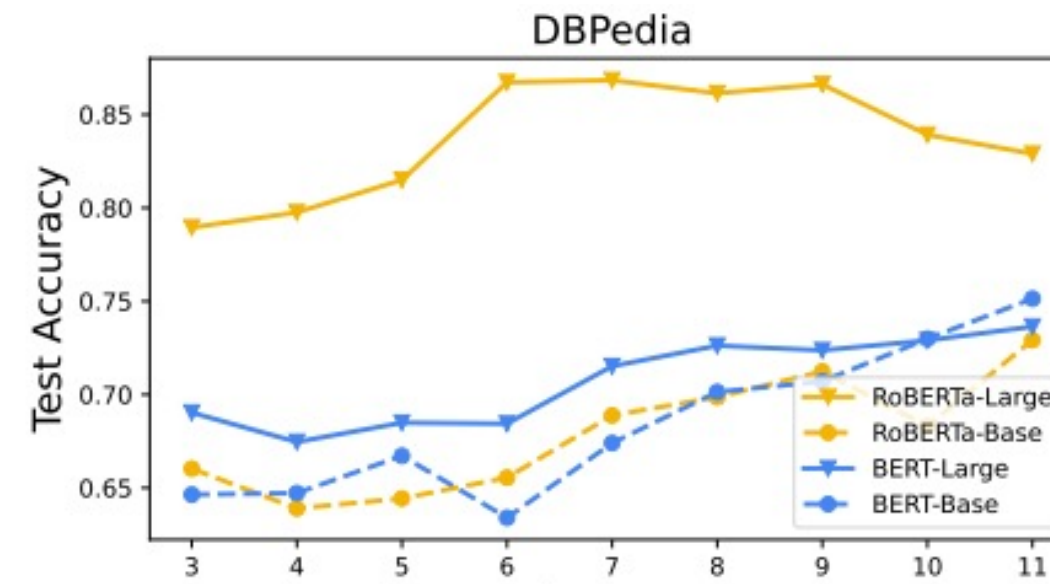
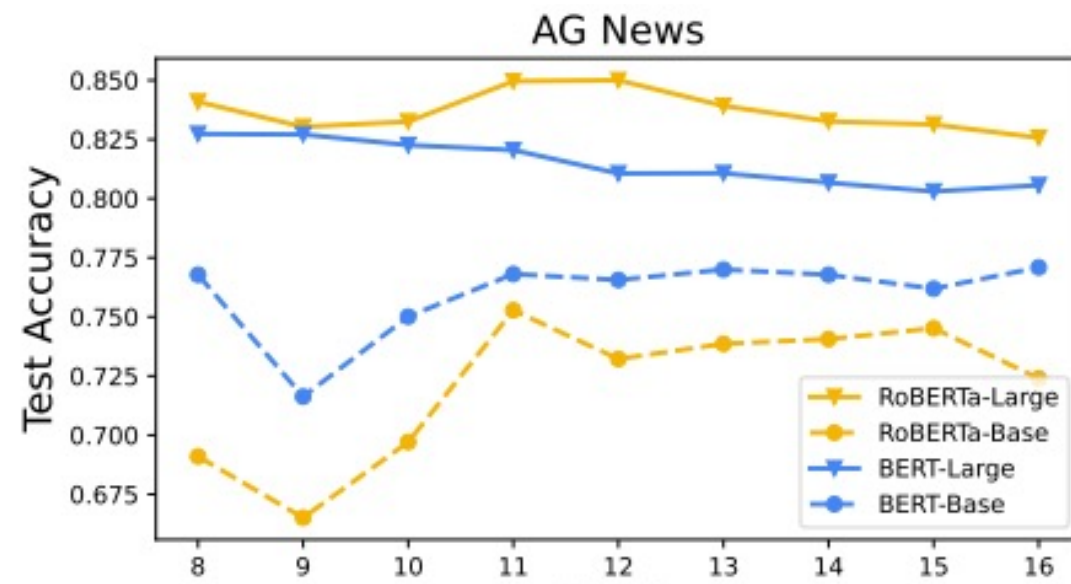


Figure 3: Test results on AG News.

# 考察③：k-近傍の単語数による性能差

- データセットによってkの値で性能差が生じる
- kを8～10くらいの設定するのが良い





## 考察④：言語モデルの違いによる性能差

### ■ 言語モデルの違いによる性能差を評価

- ・ T5-base, GPT2-base, BERT-base/large, RoBERTa-base/largeを比較
- ・ T5、GPTはモデルが最初に生成した単語を元にラベルを予測



## 考察④：言語モデルの違いによる性能差

- 言語モデルの違いで大きく性能に差が存在
  - ・ PLMの初期単語埋め込みと語彙に依存
  - ・ largeモデルはbaseモデルを上回る性能を示す

<b>Method</b>	<b>AG</b>	<b>DB</b>	<b>IM</b>	<b>AZ</b>	<b>Avg.</b>
NPPrompt-T5-base	76.8	78.3	68.5	65.3	72.2
NPPrompt-GPT2-base	81.1	78.1	83.7	85.6	82.1
NPPrompt-BERT-base	79.4	77.8	57.7	53.5	67.1
NPPrompt-BERT-large	82.7	80.9	81.6	80.8	81.5
NPPrompt-RoBERTa-base	75.3	82.8	88.7	83.9	82.7
NPPrompt-RoBERTa-large	<b>85.0</b>	<b>86.8</b>	<b>94.1</b>	<b>93.9</b>	<b>90.0</b>

## 考察④：言語モデルの違いによる性能差

- 言語モデルの違いで大きく性能に差が存在
  - ・ PLMの初期単語埋め込みと語彙に依存
  - ・ largeモデルはbaseモデルを上回る性能を示す

<b>Method</b>	<b>AG</b>	<b>DB</b>	<b>IM</b>	<b>AZ</b>	<b>Avg.</b>
NPPrompt-T5-base	76.8	78.3	68.5	65.3	72.2
NPPrompt-GPT2-base	81.1	78.1	83.7	85.6	82.1
NPPrompt-BERT-base	79.4	77.8	57.7	53.5	67.1
NPPrompt-BERT-large	82.7	80.9	81.6	80.8	81.5
NPPrompt-RoBERTa-base	75.3	82.8	88.7	83.9	82.7
NPPrompt-RoBERTa-large	<b>85.0</b>	<b>86.8</b>	<b>94.1</b>	<b>93.9</b>	<b>90.0</b>

## 考察⑤：テキスト分類以外のタスクへの適用

### ■ 質問応答へNPPromptを適用

- データセットは選択問題 ([CommonsenseQA](#))
- “x The answer is [MASK].” というプロンプト
- 比較モデル

1. GPT-J Direct / CoT

2. LaMDA (137B)

*Where on a **river** can you hold a cup upright to catch water on a sunny day?*

✓ waterfall, ✗ bridge, ✗ valley, ✗ pebble, ✗ mountain

*Where can I stand on a **river** to see water falling without getting wet?*

✗ waterfall, ✓ bridge, ✗ valley, ✗ stream, ✗ bottom

*I'm crossing the **river**, my feet are wet but my body is dry, where am I?*

✗ waterfall, ✗ bridge, ✓ valley, ✗ bank, ✗ island

[Talmor+,19]

## 考察⑤：テキスト分類以外のタスクへの適用

### ■ 質問応答へNPPromptを適用

- ・ QAでもある程度の性能
- ・ Few-shotのGPT-Jのスコアにも勝るような性能
- ・ 他のNLPタスクへも応用できる可能性

<b>Method</b>	<b>CQA Dev Set Accuracy</b>
Few-shot Direct GPT-J	20.9
Few-shot CoT GPT-J	36.6
Few-shot CoT LaMDA 137B	55.6
NPPrompt-RoBERTa-large	34.2

# まとめと今後の展望

- Zero-shotでテキスト分類などの手法の性能を向上させる手法の提案
- PLMの語彙を利用することで人手のコストや追加データの用意が不要
- 従来のZero-shotの性能を上回る結果
- Few-shotの手法にも匹敵するような性能
  
- 今後の展望
  - ・ キーワードの設計次第でラベル名が意味を持たなくても利用可能
  - ・ この手法をより大きなモデルに適用したら性能は上がりそう