

Document Alignment based on Overlapping Fixed-Length Segmentation using Optimal Transport

○王 小天¹ 宇津呂 武仁¹ 永田 昌明²

1. 筑波大学 システム情報工学研究群 自然言語処理研究室
2. NTTコミュニケーション科学基礎研究所

目次

1. 研究背景
2. 先行研究
3. 提案手法
4. データセット
5. 実験設定
6. 実験結果と分析
7. まとめ

目次

1. 研究背景

2. 先行研究

3. 提案手法

4. データセット

5. 実験設定

6. 実験結果と分析

7. まとめ

研究背景

◆ウェブクローラデータ

- 機械翻訳に関する研究では、ウェブクローリングした多言語データを利用して、
パラレルコーパスを構築することが非常に重要
- 代表的には、ParaCrawl Dataset (Bañón et al., ACL 2020)、JParaCrawl Dataset (Morishita et al., LREC 2022)など
- ウェブクローラデータの作成手順：クローリング、テキスト抽出、
二言語間の文書対応付け、二言語間の文対応付け、文ペアのフィルタリング

研究背景

◆二言語間の文書対応付け

■定義

異なる文書集合の中でお互いに対訳になっている文書を結びつけること

■目的

ウェブクローラから取得した多言語の文書に対して、対訳コーパスを作るために、文書の二言語間対応関係を作る

目次

1. 研究背景

2. 先行研究

3. 提案手法

4. データセット

5. 実験設定

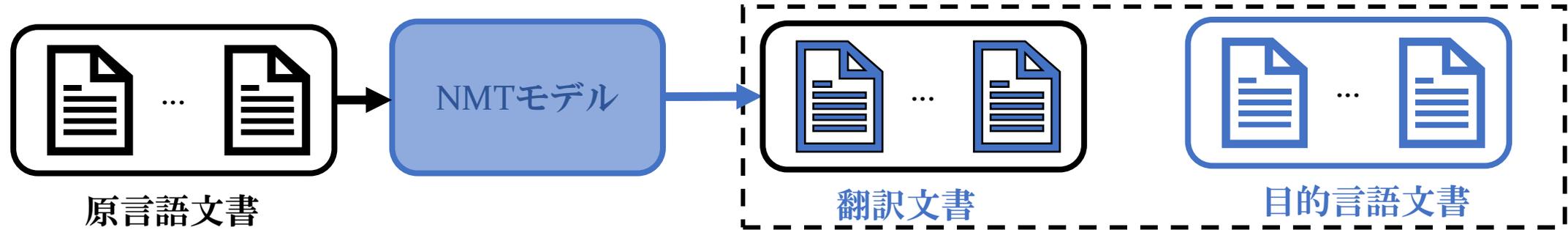
6. 実験結果と分析

7. まとめ

先行研究 | 機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

◆ 機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

■ 手順

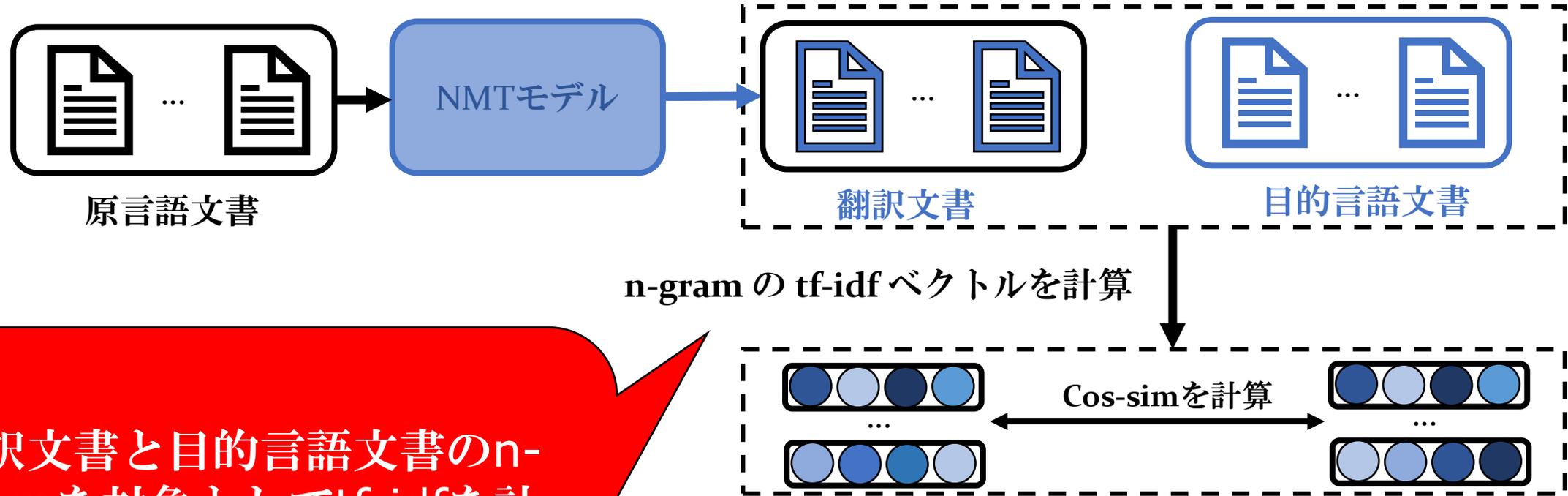


原言語文書の文単位分割された文を対象として、目的言語に機械翻訳する

先行研究 | 機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

◆ 機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

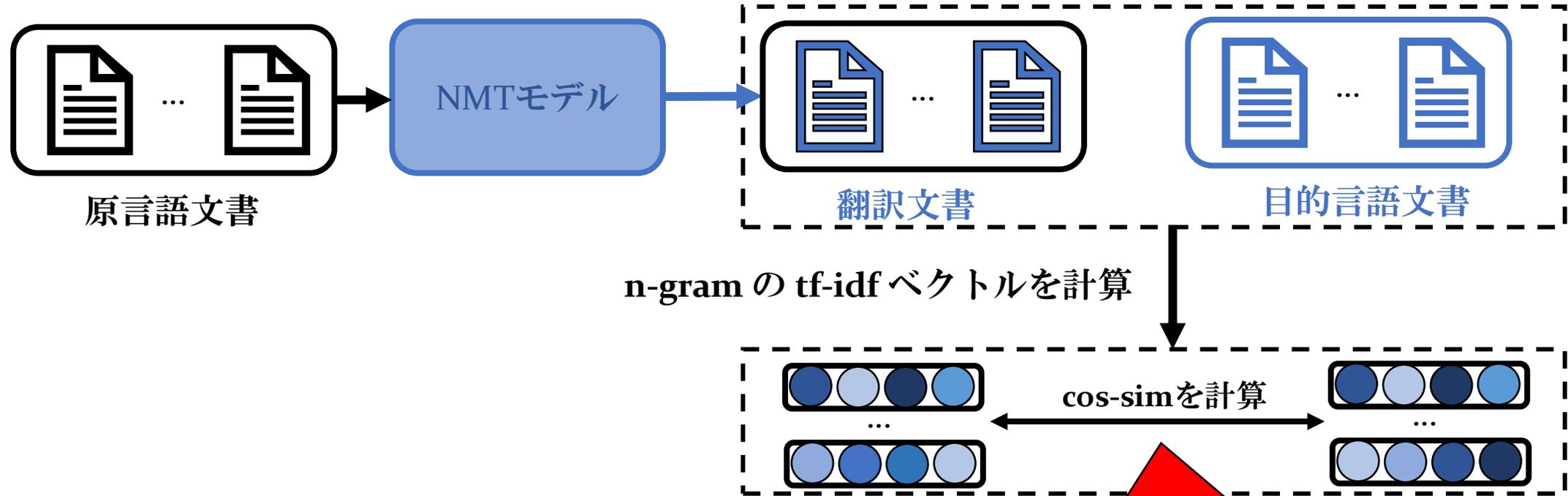
■ 手順



翻訳文書と目的言語文書のn-gramを対象としてtf-idfを計算し、各文書に対してtf-idfベクトルを生成する

◆機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

■手順



翻訳文書と目的言語文書のtf-idfベクトル間の cosine類似度に基づいて、対応関係を作る

◆ 機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

■ 対応関係の作成ルール (1-1ルール)

- 二言語間文書対の全候補を類似度の高い順に並べ替え、対応先が決まった文書対は削除
- 各文書が最大で一つの対応文書しか持たないようにする

並べ替え後の二言語間文書対

日本語文書-342

英語文書-101

日本語文書-12

英語文書-101

日本語文書-5

英語文書-42

日本語文書-342

英語文書-9

◆ 機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

■ 対応関係の作成ルール (1-1ルール)

- 二言語間文書対の全候補を類似度の高い順に並べ替え、対応先が決まった文書対は削除
- 各文書が最大で一つの対応文書しか持たないようにする

並べ替え後の二言語間文書対

日本語文書-342

英語文書-101

日本語文書-12

~~英語文書-101~~

日本語文書-5

英語文書-42

~~日本語文書-342~~

英語文書-9

対応先が決まった文書対は削除

◆ 機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

■ 対応関係の作成ルール (1-1ルール)

- 二言語間文書対の全候補を類似度の高い順に並べ替え、対応先が決まった文書対は削除
- 各文書が最大で一つの対応文書しか持たないようにする

並べ替え後の二言語間文書対

日本語文書-342

英語文書-101

日本語文書-12

~~英語文書-101~~

日本語文書-5

英語文書-42

~~日本語文書-342~~

英語文書-9

最終結果

日本語文書-342

英語文書-101

日本語文書-5

英語文書-42

◆ Sentence Embeddingに基づく二言語間の文書対応付け (Mean Pooling)

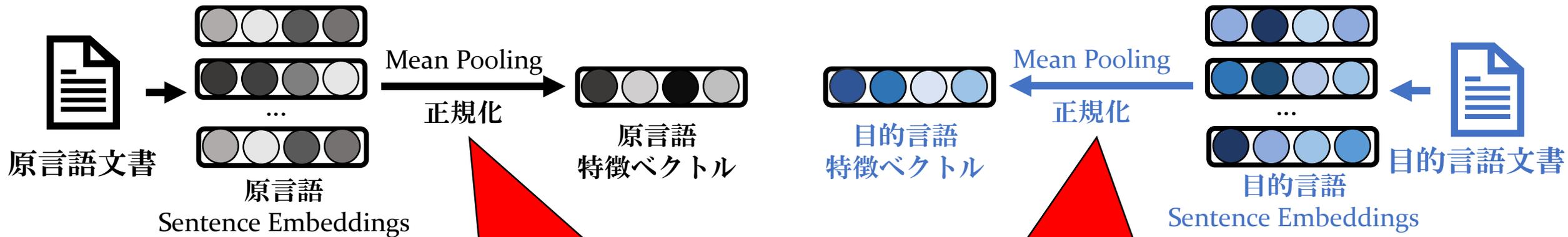
■ 手順



原言語文書と目的言語文書の分割された文を対象として、事前訓練済み多言語Sentence Embeddingモデルにより、それぞれに文埋め込みに変換

◆ Sentence Embeddingに基づく二言語間の文書対応付け (Mean Pooling)

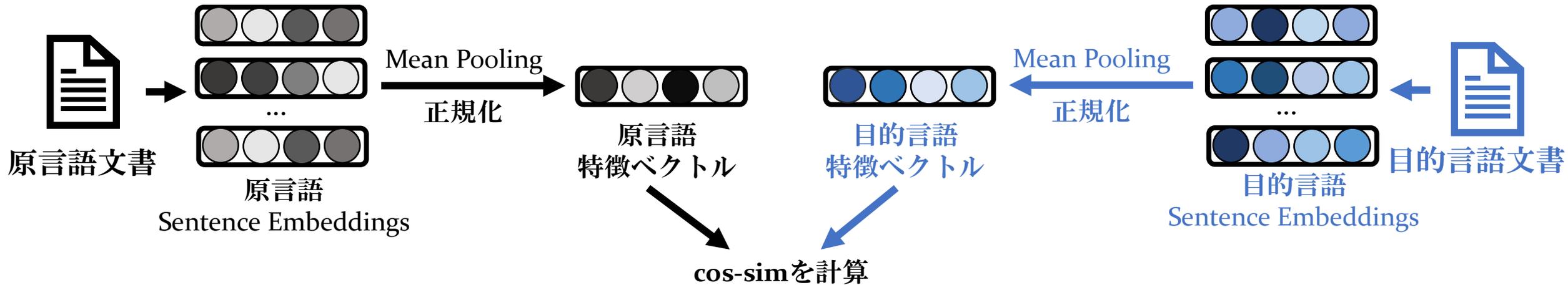
■ 手順



各文書の系列ベクトルをMean Poolingし、次に正規化して、文書の特徴ベクトルを取得する

◆ Sentence Embeddingに基づく二言語間の文書対応付け (Mean Pooling)

■ 手順



原言語文書と目的言語文書間の特徴ベクトルのcosine類似度に基づいて、対応関係を作る (1-1ルールに従う)

◆ Sentence Embeddingに基づく二言語間の文書対応付け (Optimal Transport)

(Clark et al., ACL 2019; El-Kishky and Guzmán, AACL 2020)

■ 手順

1. 原言語文書と目的言語文書に対して、事前訓練済み多言語Sentence Embeddingモデルで、それぞれの系列ベクトルを生成する
2. Sentence-Levelの文書間の最適輸送距離 (Optimal Transport Distance)、Sentence Mover's Distance (SMD)を計算する

(最適輸送距離 (OT)の計算方式は次スライド)

◆ Sentence Embeddingに基づく二言語間の文書対応付け (Optimal Transport)

(Clark et al., ACL 2019; El-Kishky and Guzmán, AACL 2020)

■ OTの計算方式

- あるドキュメントDoc A, Doc B
- Doc A の s_i の重みは $d_{A,i}$ とする
- Doc B の s'_j の重みは $d_{B,j}$ とする
- $c(i,j)$ は s_i の Embedding と s'_j の Embedding の距離
- T_{ij} の i 行目の和は $d_{A,i}$ 、 j 列目の和は $d_{B,j}$ が条件として、最小の $\sum_{i,j=1}^n T_{ij}c(i,j)$ を見つけるのが目的
- 見つかった minimum $\sum_{i,j=1}^n T_{ij}c(i,j)$ が最終の文間の最適輸送距離
- 重み付け方式の候補は右記

目標関数 :

$$\min_{T \geq 0} \sum_{i=1}^n \sum_{j=1}^m T_{ij} c(i,j)$$

Subject to:

$$\sum_{j=1}^m T_{ij} = d_{A,i}$$
$$\sum_{i=1}^n T_{ij} = d_{B,j}$$

重み付け方式

- 文頻度による重み
- 文長による重み
- IDFによる重み
- 文長、IDFによる重み

目次

1. 研究背景
2. 先行研究
- 3. 提案手法**
4. データセット
5. 実験設定
6. 実験結果と分析
7. まとめ

◆提案手法：Sentence Embeddingに基づく二言語間の文書対応付け

1. 2特徴ベクトルによる文書対応

- Mean Poolingベクトルを利用するだけでなく、文書の系列ベクトルの中の第1文ベクトルも使用
- 文書間の類似度を計算する際に、第1文ベクトルとMean Poolingベクトルの間の類似度をそれぞれに計算して、それらに重みを与えて合計する

2. Overlapping Fixed-Length Segmentation (OFLS)による文書対応

- 文単位分割ではなく、固定長 L の Sliding Window で、ウェブクロールされた文書を重複ありの Segments に分割する
- OFLSをもとに、言語の影響を考慮して、原言語文書と目的言語文書を異なる固定長(言語間で一定の比率)で分割する (Language-Pair Dependent OFLS, LD-OFLS)
- 従来手法・提案手法において、「文単位分割」と「OFLS」の性能を比較する

◆提案手法：Sentence Embeddingに基づく二言語間の文書対応付け

1. 2特徴ベクトルによる文書対応

- ▶ Mean Poolingベクトルを利用するだけでなく、文書の系列ベクトルの中の第1文ベクトルも使用
- ▶ 文書間の類似度を計算する際に、第1文ベクトルとMean Poolingベクトルの間の類似度をそれぞれに計算して、それらに重みを与えて合計する

2. Overlapping Fixed-Length Segmentation (OFLS)による文書対応

- ▶ 文単位分割ではなく、固定長 L の Sliding Window で、ウェブクロールされた文書を重複ありの Segments に分割する
- ▶ OFLSをもとに、言語の影響を考慮して、原言語文書と目的言語文書を異なる固定長(言語間で一定の比率)で分割する (Language-Pair Dependent OFLS, LD-OFLS)
- ▶ 従来手法・提案手法において、「文単位分割」と「OFLS」の性能を比較する

◆Sentence Embeddingに基づく二言語間の文書対応付け (2 特徴ベクトル)

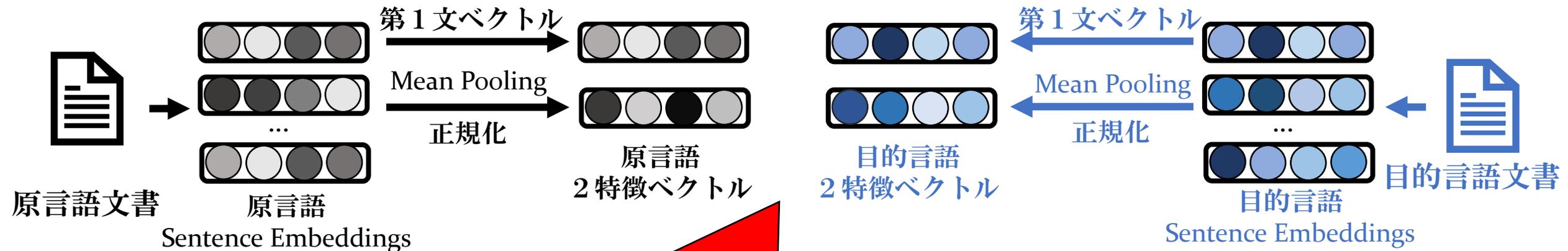
■手順



原言語文書と目的言語文書に対して、
それぞれの系列ベクトルを生成する

◆Sentence Embeddingに基づく二言語間の文書対応付け (2 特徴ベクトル)

■手順

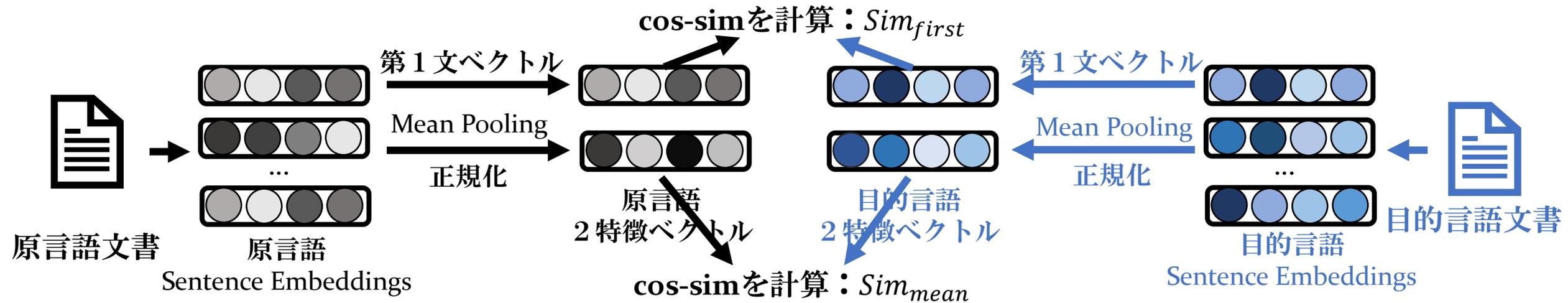


各文書の系列ベクトルのMean Poolingベクトルと第1文ベクトルを取得する

◆Sentence Embeddingに基づく二言語間の文書対応付け (2 特徴ベクトル)

■手順

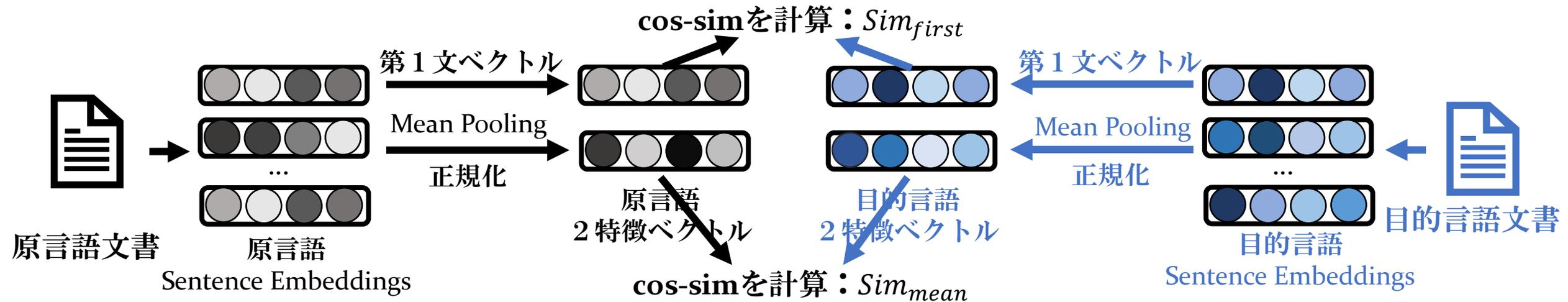
第1文ベクトル間のcosine類似度を計算



Mean Poolingベクトル間のcosine類似度を計算

◆Sentence Embeddingに基づく二言語間の文書対応付け (2 特徴ベクトル)

■手順



類似度の重み付き和を計算: $Sim_{final} = \lambda Sim_{first} + (1 - \lambda) Sim_{mean}$

◆提案手法：Sentence Embeddingに基づく二言語間の文書対応付け

1. 2特徴ベクトルによる文書対応

- ▶ Mean Poolingベクトルを利用するだけでなく、文書の系列ベクトルの中の第1文ベクトルも使用
- ▶ 文書間の類似度を計算する際に、第1文ベクトルとMean Poolingベクトルの間の類似度をそれぞれに計算して、それらに重みを与えて合計する

2. Overlapping Fixed-Length Segmentation (OFLS)による文書対応

- ▶ 文単位分割ではなく、固定長 L の Sliding Window で、ウェブクロールされた文書を重複ありの Segments に分割する
- ▶ OFLSをもとに、言語の影響を考慮して、原言語文書と目的言語文書を異なる固定長(言語間で一定の比率)で分割する (Language-Pair Dependent OFLS, LD-OFLS)
- ▶ 従来手法・提案手法において、「文単位分割」と「OFLS」の性能を比較する

◆ Overlapping Fixed-Length Segmentation (OFLS)による文書対応

■ 定義

- ・ 固定長 L の Sliding Window を利用して Segments に分割
- ・ 隣接する Segmentsの間では、指定された割合が重複

■ OFLSの例

今日は晴れていて、気温は25度ほどで、風も穏やかです。公園で友達とピクニックを楽しむ予定です。お弁当にはおにぎり、サンドイッチ、フルーツ、そして冷たいジュースが入っています。みんなで楽しい時間を過ごし、笑顔が絶えないことを願っています。

$L = 20$

重複割合: 0.5

Seg1: 今日は晴れていて、気温は25度ほどで、風も穏やかです。公園で
Seg2: ほどで、風も穏やかです。公園で友達とピクニックを楽しむ予定です。お弁当
Seg3: 友達とピクニックを楽しむ予定です。お弁当にはおにぎり、サンドイッチ、フルーツ、そして冷たい
Seg4: にはおにぎり、サンドイッチ、フルーツ、そして冷たいジュースが入っています。
みんなで楽しい時間
Seg5: ジュースが入っています。みんなで楽しい時間を過ごし、笑顔が絶えないことを願っています。

◆ Language-Pair Dependent OFLS (LD-OFLS)による文書対応

■ 考慮点

一般的に、異なる言語で同じ意味を表現するために必要なトークン数は異なる

例：「私は犬が好きだ」と「I like dogs」

日本語：6トークン、英語：3トークン

⇒ 固定長分割において、原言語文書・目的言語文書で異なる固定長を使用

■ 手法

- ρ (language-pair dependent proportion) : 言語間の文書長(トークン数)の比率
- 目的言語文書を分割する固定長 = 原言語文書を分割する固定長 * ρ

◆Language-Pair Dependent OFLS (LD-OFLS)による二言語間文書対応付け

■具体例

The expected annual power generation will be 87,000,000kWh, corresponding to the annual electricity consumption of 30,000 ordinary houses.

「setu4993/LaBSE」でtokenize

The expected annual power generation will be 87,000,000 ##kw ##h ...

総トークン数 : 28

総トークン数 : 42

年間予想発電量は8,700万キロワット時で、一般家庭...

「setu4993/LaBSE」でtokenize

年間予想発電量は8,700万キロワット時で、一般家庭約3万世帯分の年間消費電力量に相当します。

目次

1. 研究背景
2. 先行研究
3. 提案手法
- 4. データセット**
5. 実験設定
6. 実験結果と分析
7. まとめ

データセット

◆テストデータ

■データセットの構築

- 4つのweb-domain **丸紅**、**西新宿**、**楽天**、**NTT CS研究**に対して、各ドメインからランダムに100件ずつ日本語文書をサンプリングし、4つの手法で対訳となる英語文書を検索する
 - 機械翻訳 + BM25
 - 機械翻訳 + tf-idf
 - URL Matching
 - CC-Aligned (El-Kishky et al., EMNLP 2020)
- 人手で参照用日英対訳文書対を作成
- 検索結果の英語文書集合を、対応目的言語(英語)文書の候補集合として利用

4つのweb-domain

- www.marubeni.com
- nishishinjuku.co.jp
- corp.rakuten.co.jp___global.rakuten.com
- www.kecl.ntt.co.jp

	丸紅	西新宿	楽天	NTT CS研究	ALL
対象： 対訳あり日本語文書数	73	16	75	68	232
検索結果の 候補英語文書数	251	42	319	319	931

データセット

◆テストデータ

■データセットの情報

- 丸紅、西新宿、楽天、NTT、これらの4つのドメイン、およびそれらの混合データ (4つのドメインの混合データ)
- 各文書の長さはHugging Faceの「setu4993/LaBSE」でtokenizeした後のトークン数
- ρ = 正解の英語文書の平均トークン数/対訳あり日本語文書の平均トークン数 (Ja-En方向)

	丸紅	西新宿	楽天	NTT CS研究	ALL
対象： 対訳あり日本語文書数	73	16	75	68	232
検索結果の 候補英語文書数	251	42	319	319	931
対訳あり日本語文書の 平均トークン数	2447.36	340.56	3541.55	726.37	2151.35
正解の英語文書の 平均トークン数	1598.51	217.69	2174.49	441.88	1350.47
ρ	0.65	0.64	0.61	0.61	0.63

目次

1. 研究背景
2. 先行研究
3. 提案手法
4. データセット
- 5. 実験設定**
6. 実験結果と分析
7. まとめ

実験設定

◆実験設定

■ Baseline: 機械翻訳に基づく二言語間の文書対応付け (Buck and Koehn, WMT2016)

- 翻訳機 : JParaCrawl v3.0 ja-en model
- 翻訳方向 : 日本語から英語
- tf-idf で文書対応 : Bitextorの「Docalign」ツールを使用 (bigramベースで実行)

■ Sentence Embeddingに基づく二言語間の文書対応付け

- 文埋め込みモデルとtokenizer : Hugging Face の「setu4993/LaBSE」
- 対訳候補の検索 : OTの計算において計算時間を要する ⇒
Mean Pooling ベクトルのcosine類似度による「Faiss」で20個の対訳候補を検索
- OTの重み付け方式 : 文頻度

■ GPU : NVIDIA RTX A6000一枚 を使用

実験設定

◆ 実験設定

■ 検索方向と範囲

- 方向：日英方向 (日本語文書に対応する英語文書を検索)
- 範囲：4つのweb-domainを混合して、全範囲で検索する

■ 評価指標

- 適合率 (Precision)
- 再現率 (Recall)
- F1_score : $\frac{2*Precision*Recall}{Precision + Recall}$

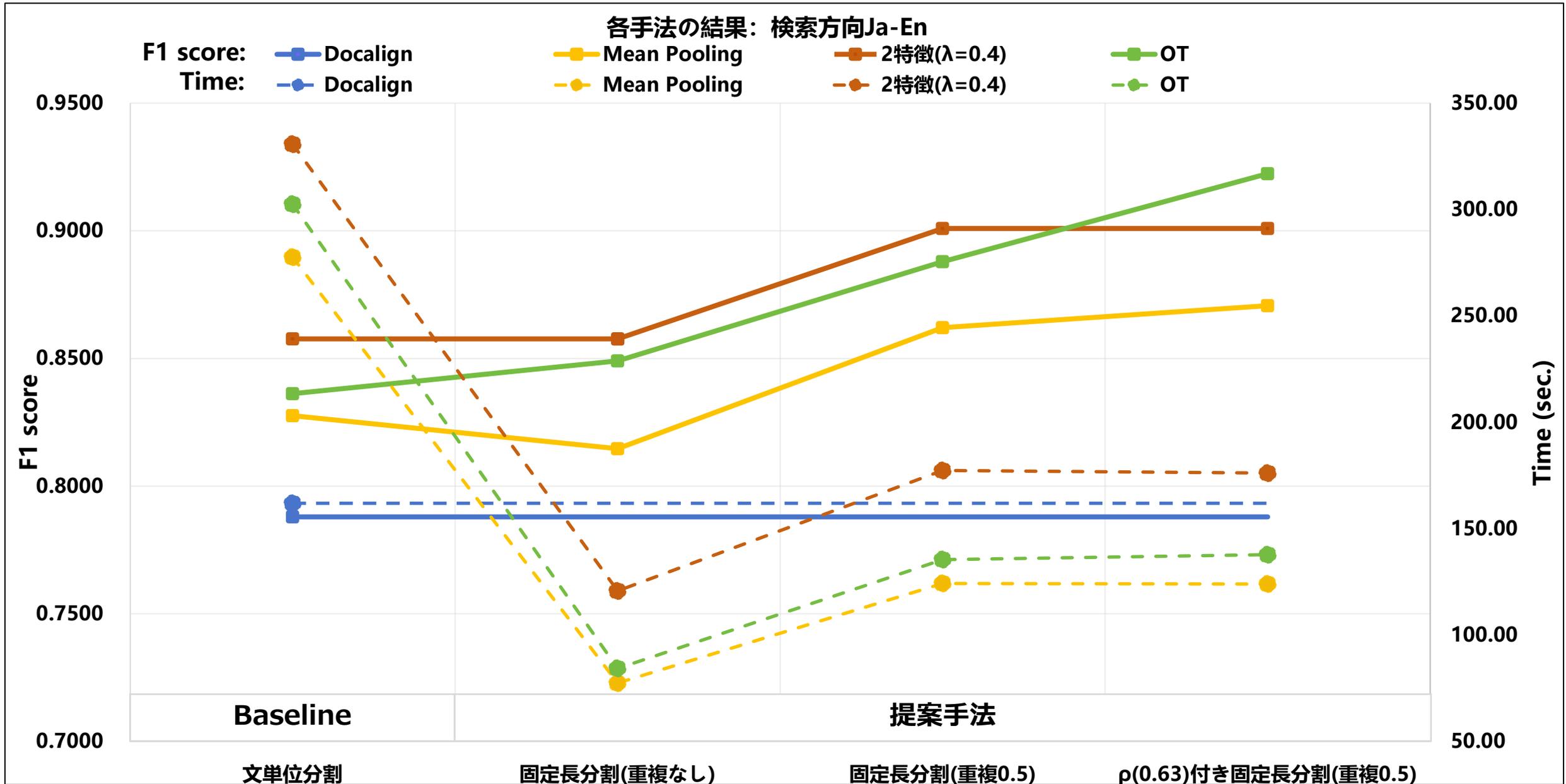
■ 計算時間の評価

- 機械翻訳に基づく二言語間の文書対応付け
 - 機械翻訳時間
 - 「Docalign」による二言語間文書対応付け時間
- Sentence Embeddingに基づく二言語間の文書対応付け
 - 系列ベクトルの生成時間、またはそれに加えて特徴ベクトルの生成時間
 - 二言語間文書対応付け時間 (類似度の計算と候補文書の検索時間が含まれる)

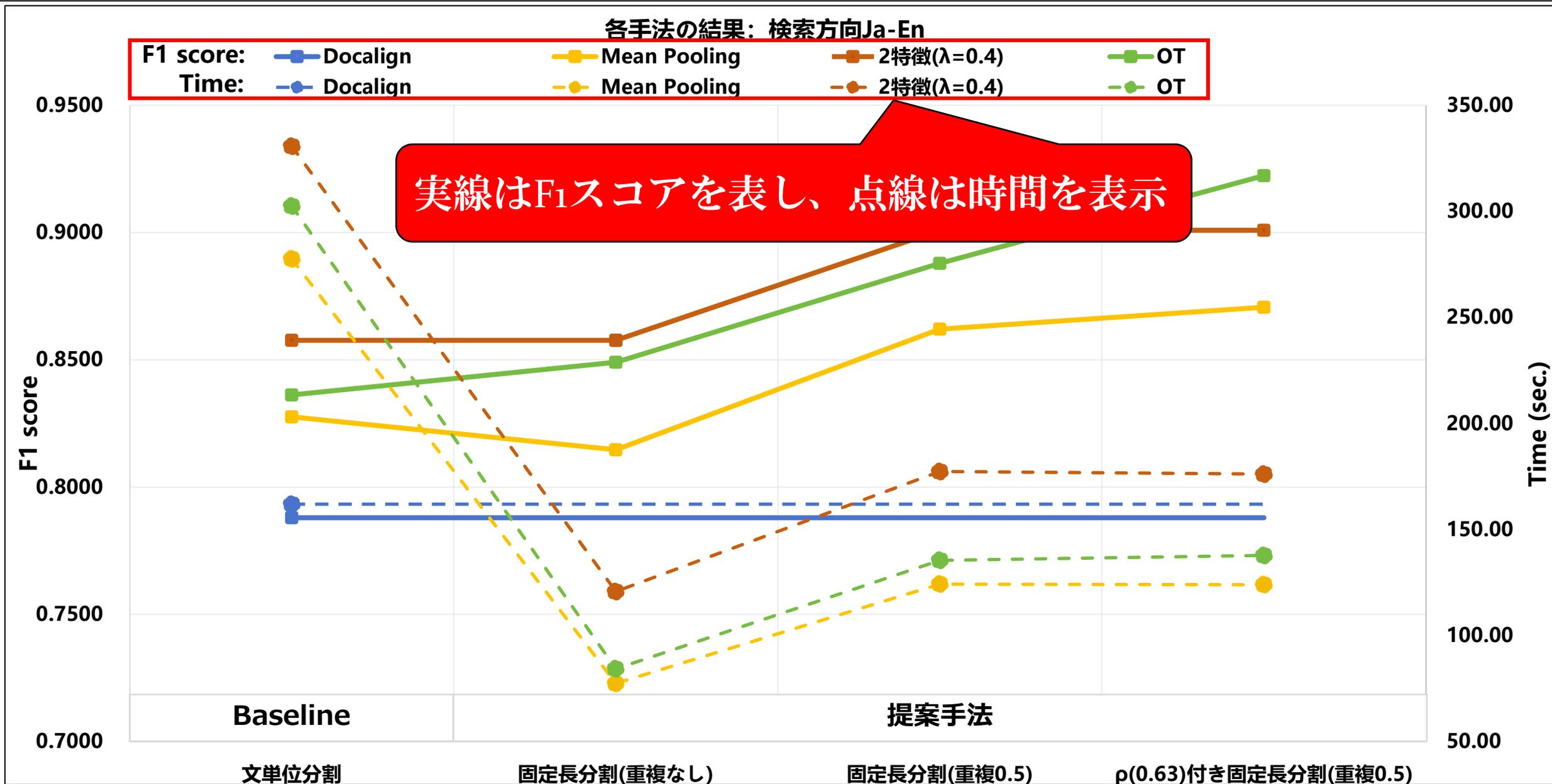
目次

1. 研究背景
2. 先行研究
3. 提案手法
4. データセット
5. 実験設定
- 6. 実験結果と分析**
7. まとめ

実験結果

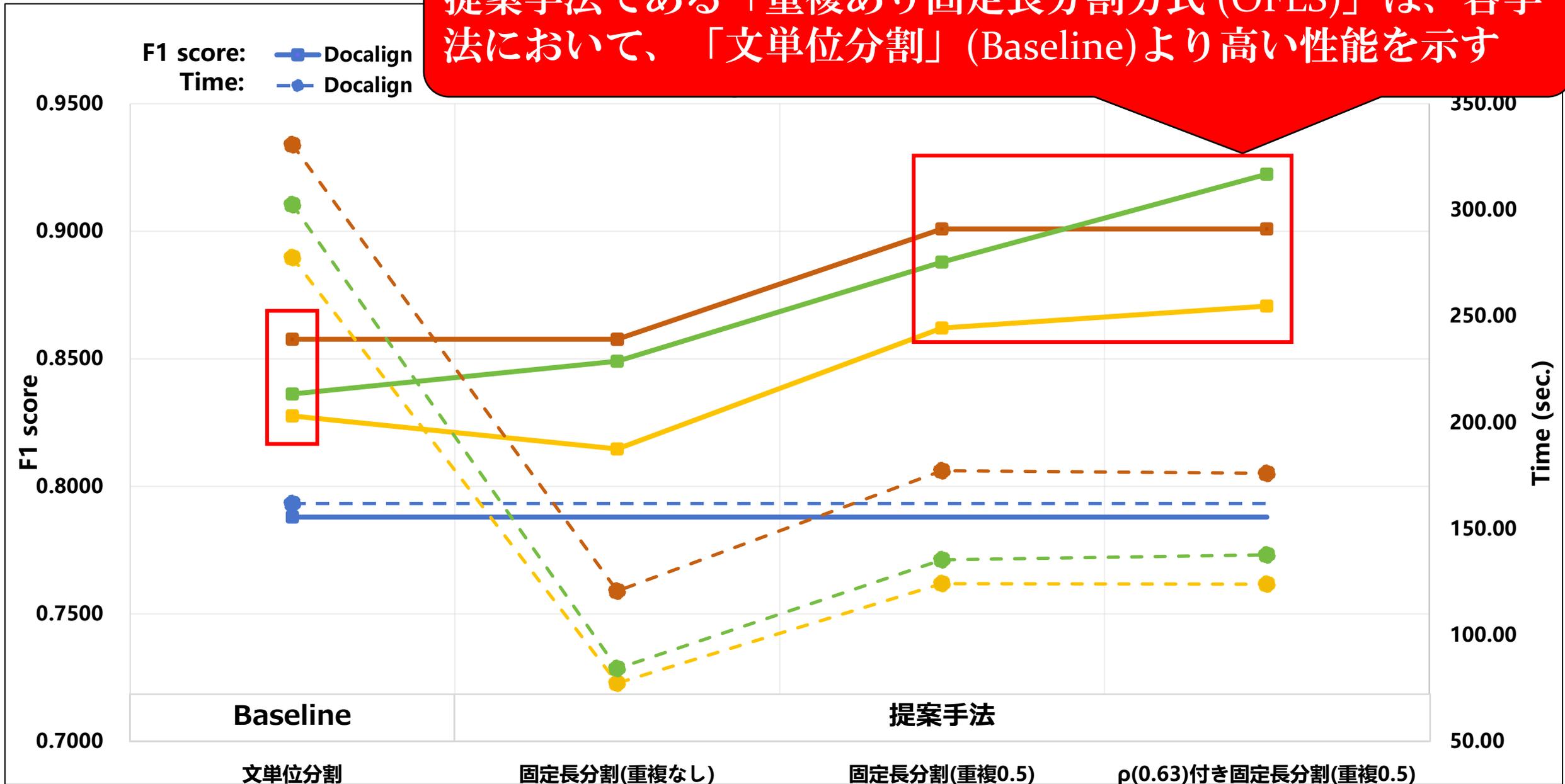


実験結果



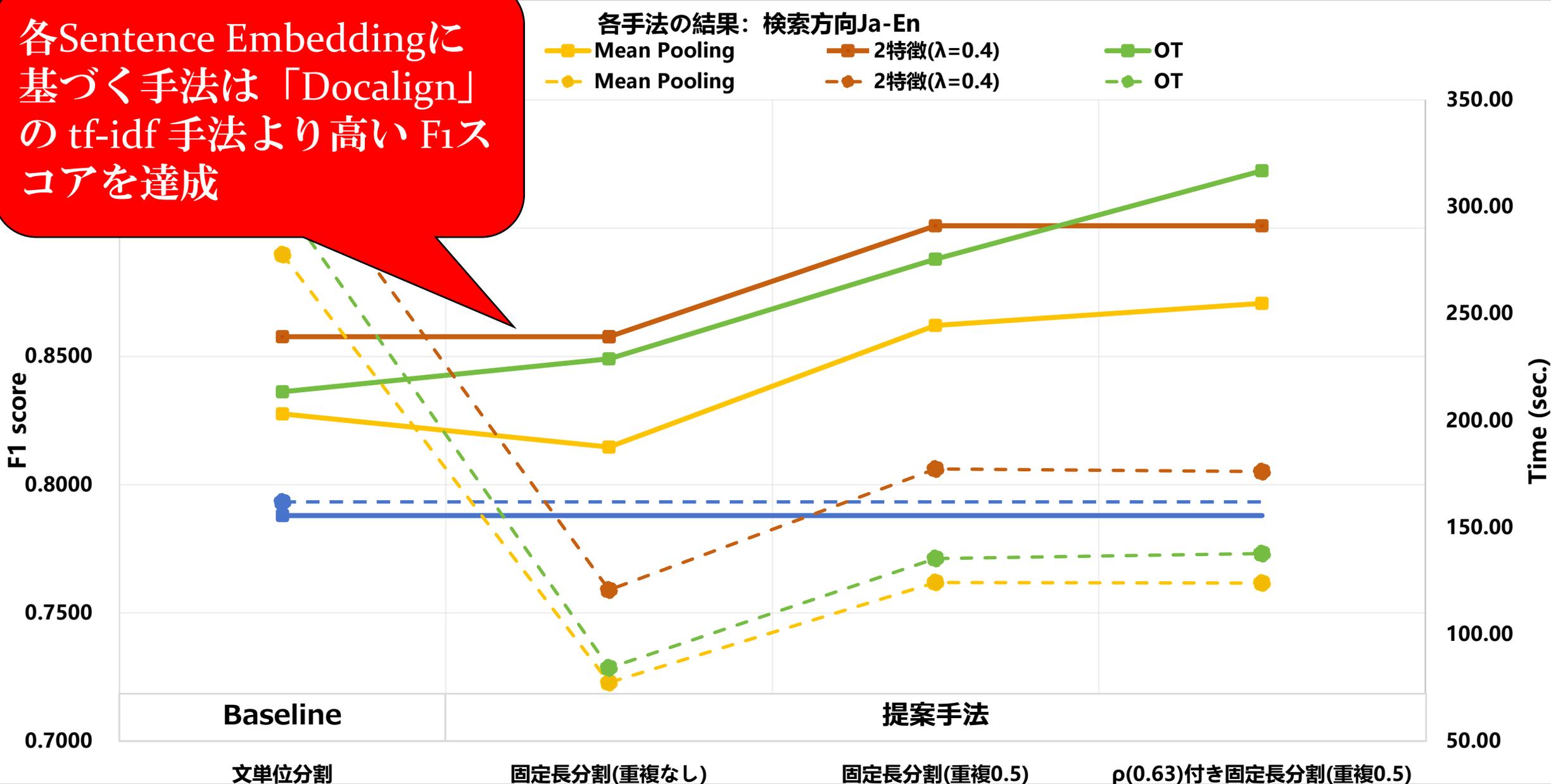
実験結果

提案手法である「重複あり固定長分割方式 (OFLS)」は、各手法において、「文単位分割」(Baseline)より高い性能を示す

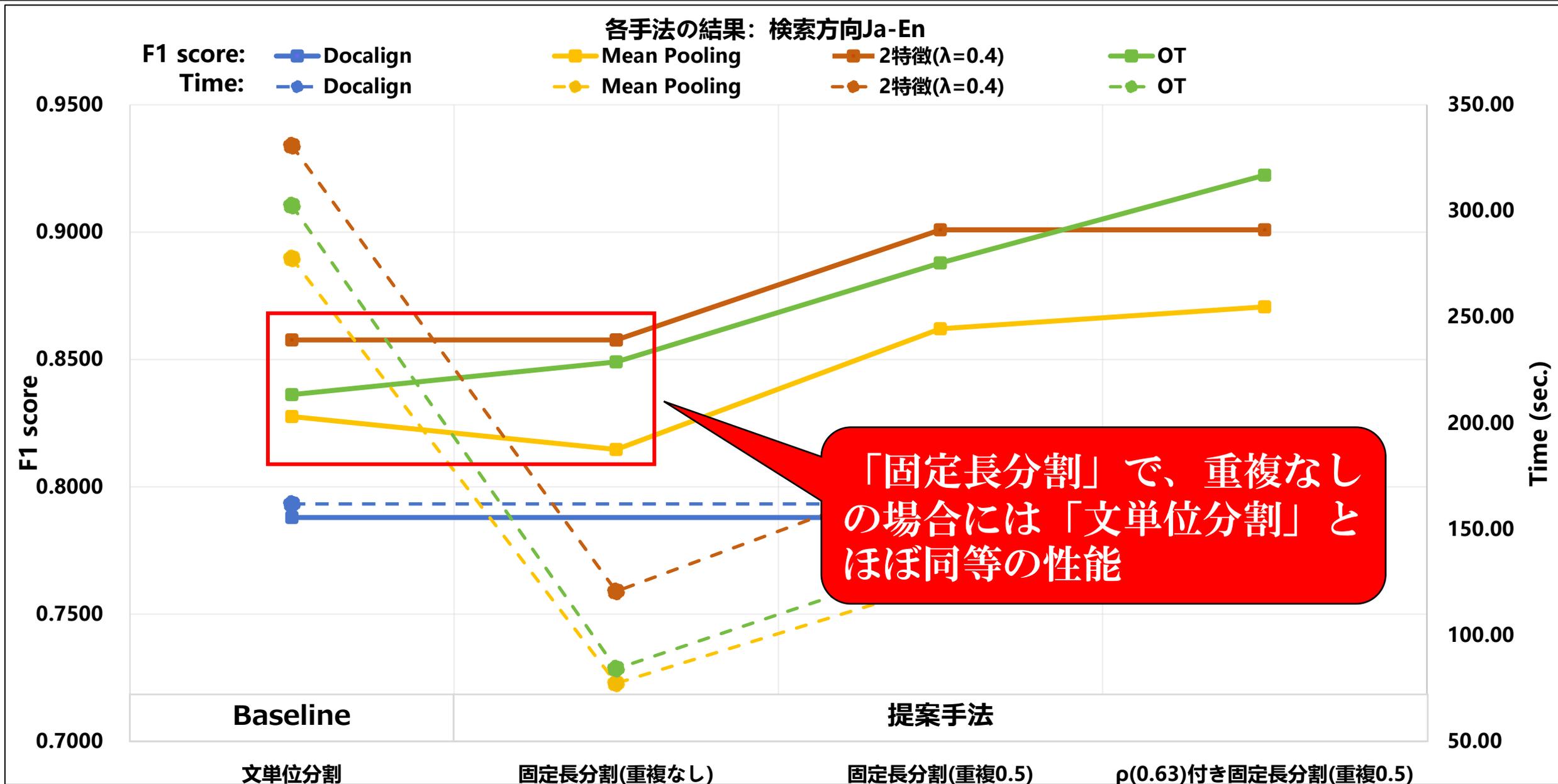


実験結果

各Sentence Embeddingに基づく手法は「Docalign」のtf-idf手法より高いF1スコアを達成

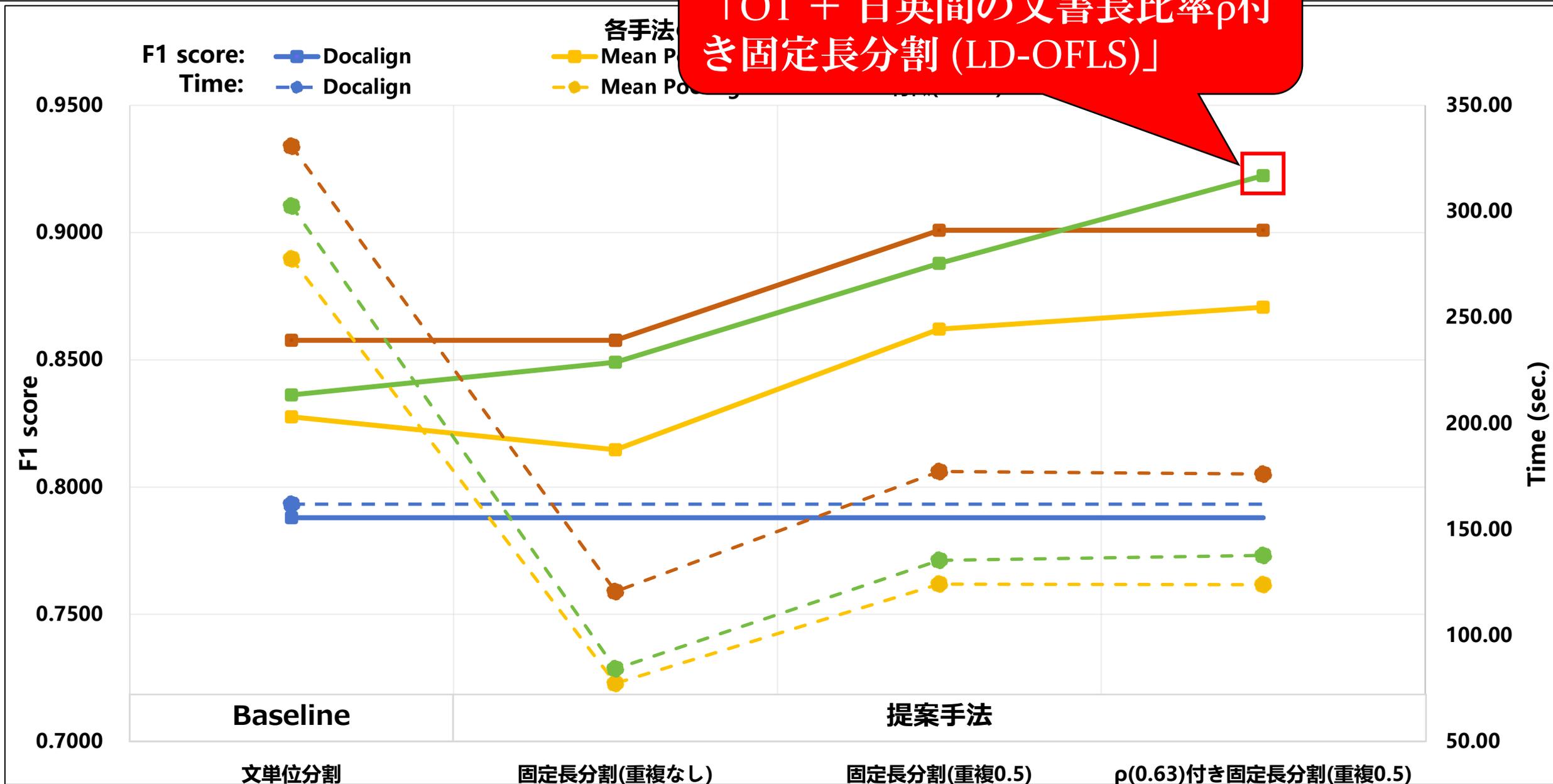


実験結果

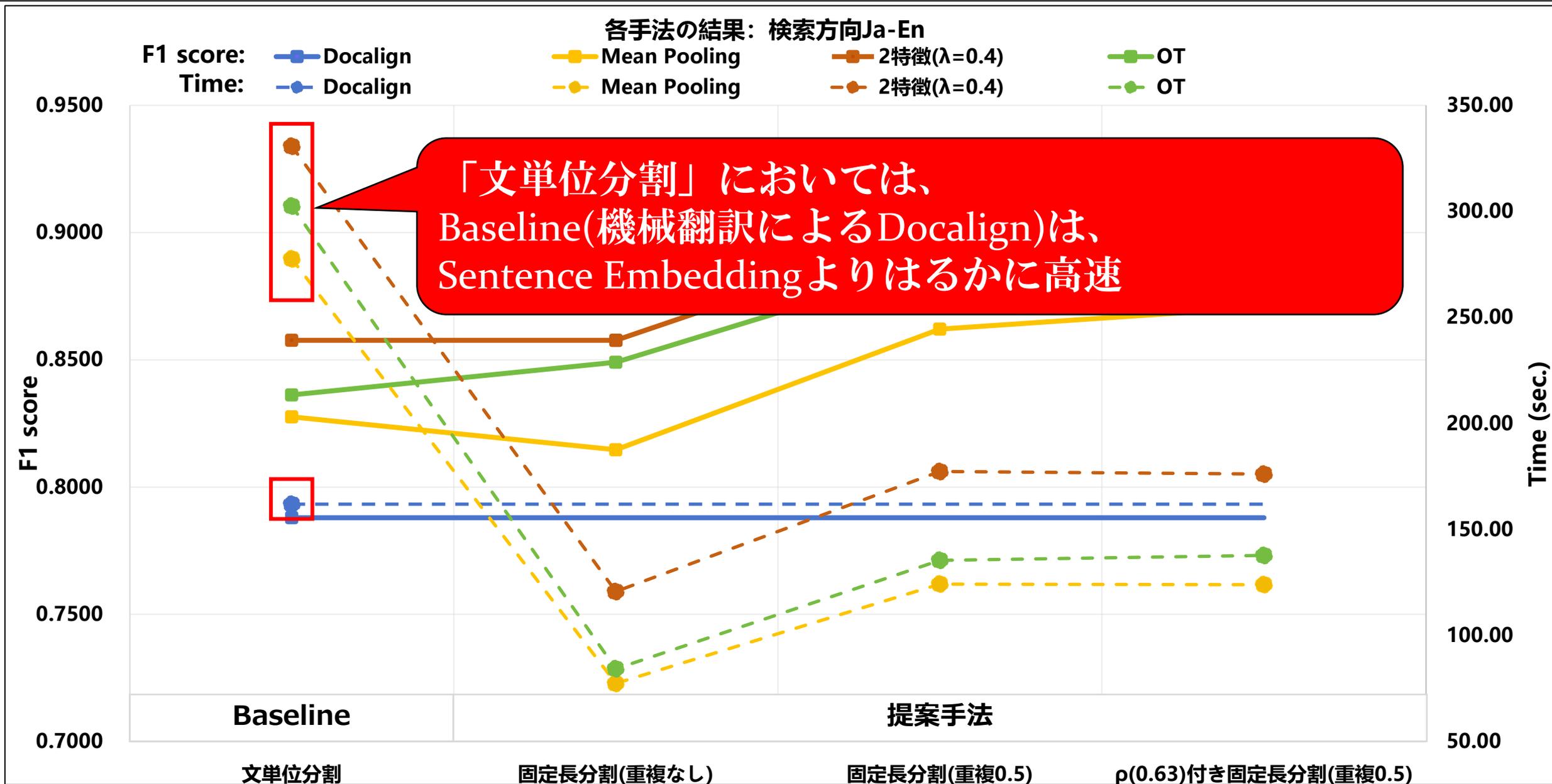


実験結果

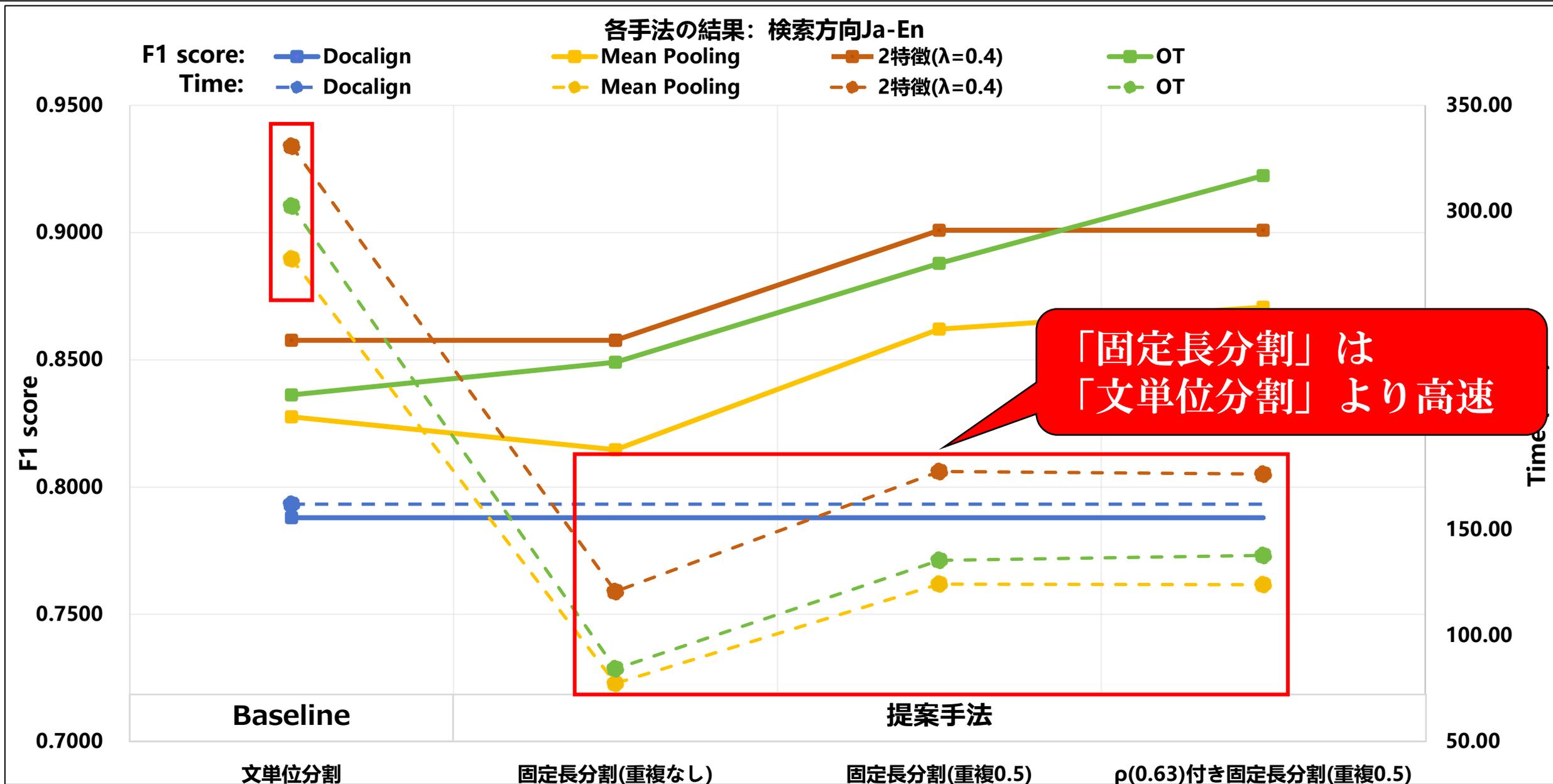
F1スコア最大は、
「OT + 日英間の文書長比率 ρ 付き
固定長分割 (LD-OFLS)」



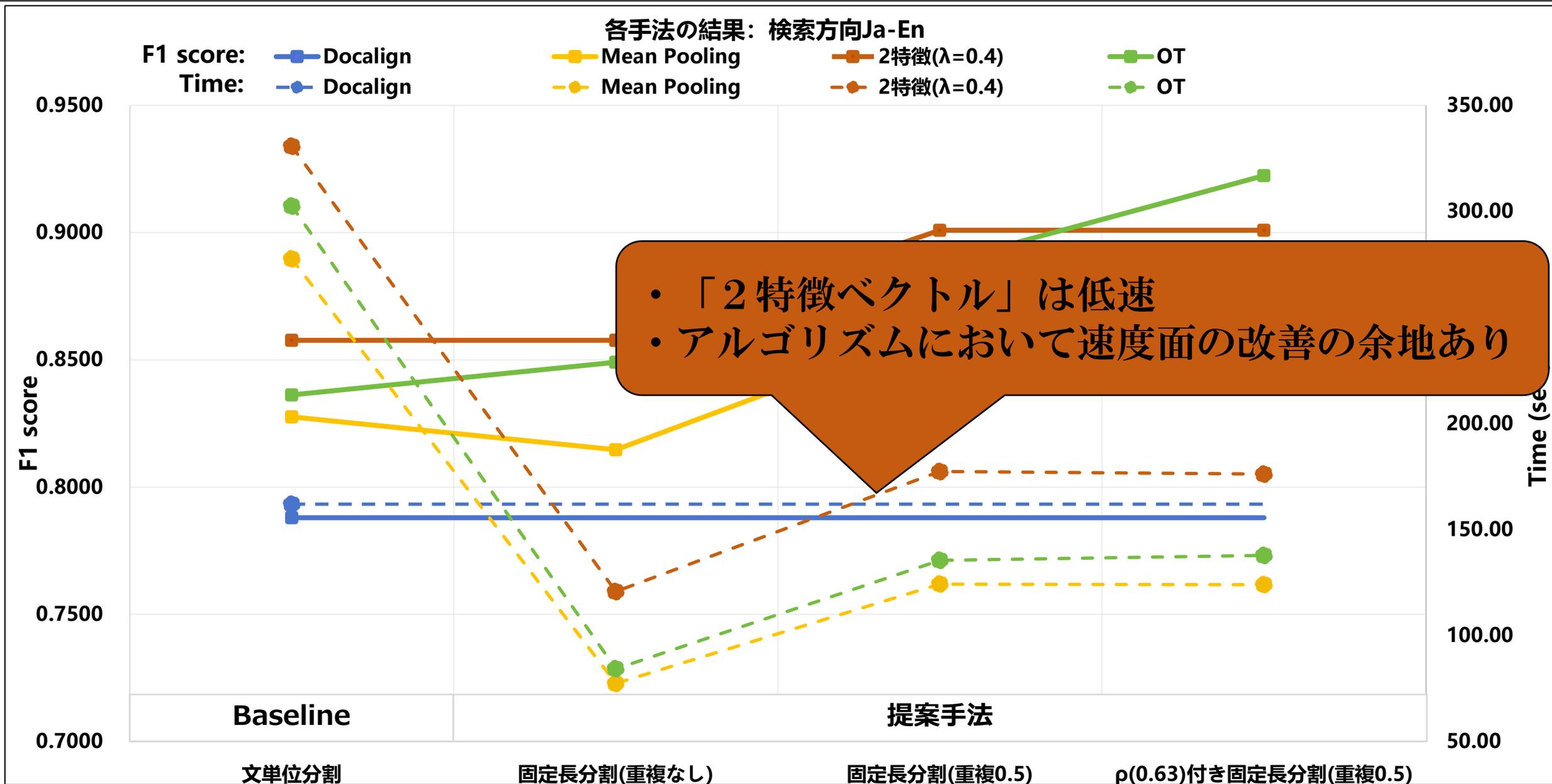
実験結果



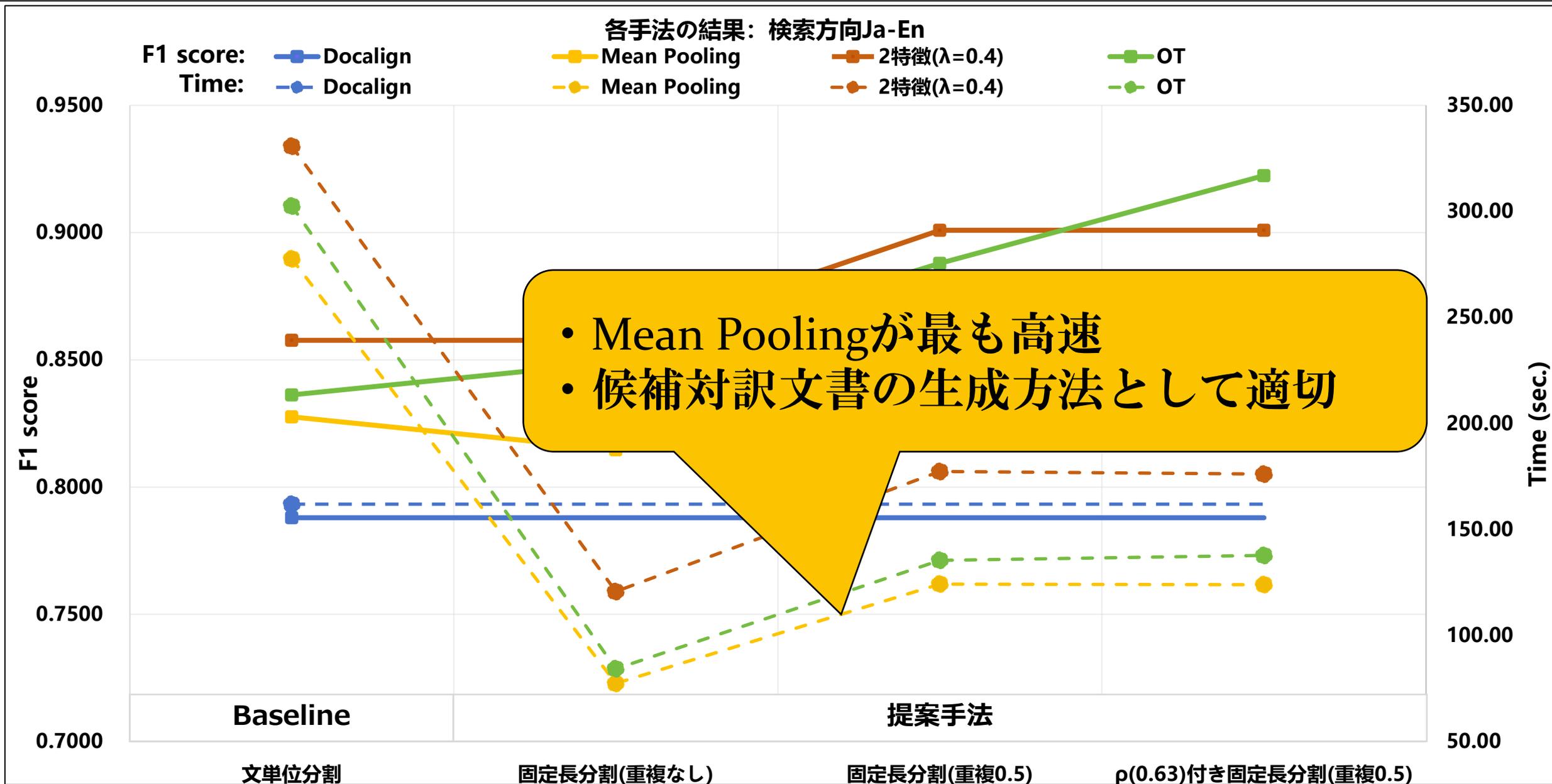
実験結果



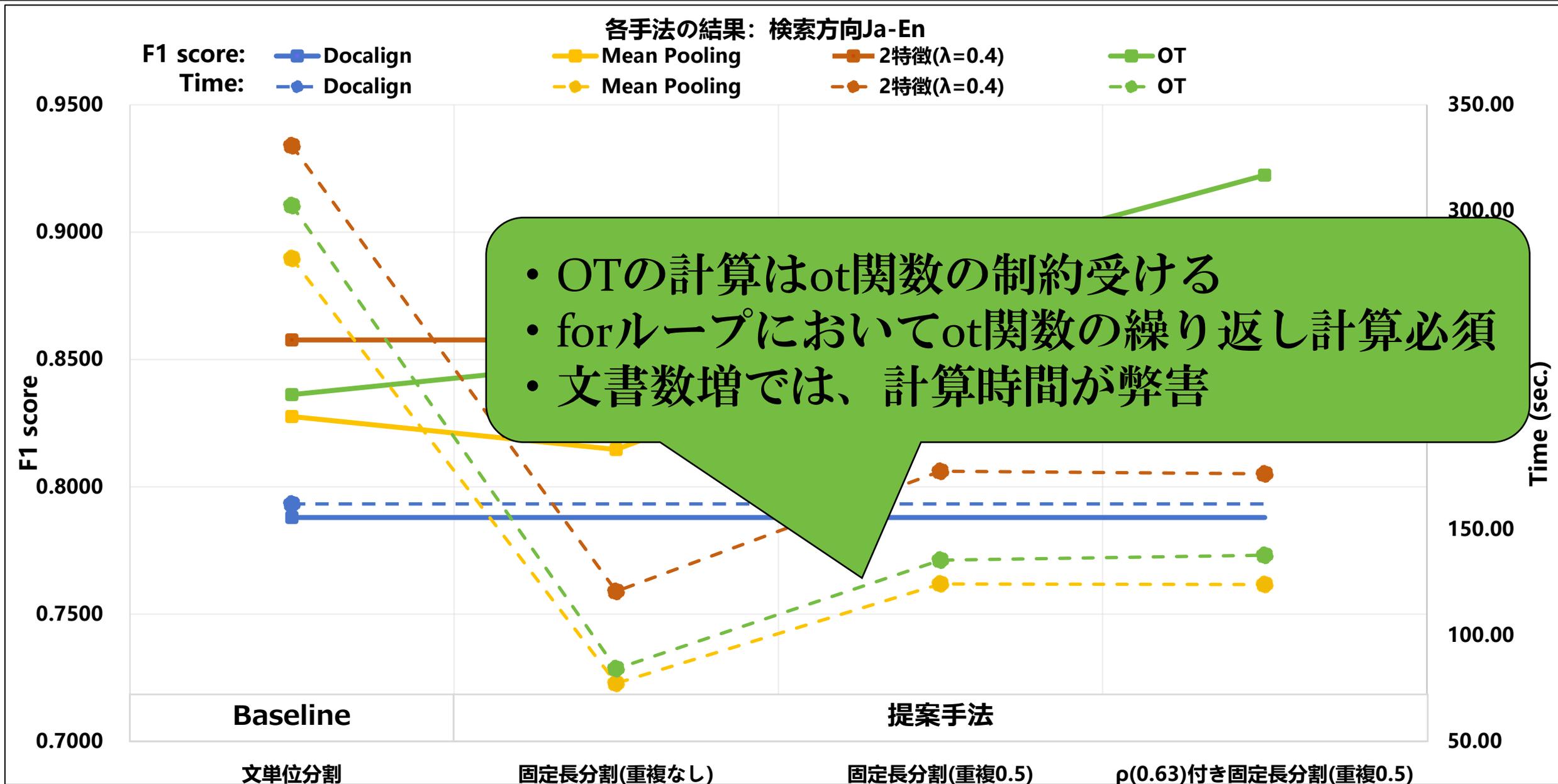
実験結果



実験結果



実験結果



目次

1. 研究背景
2. 先行研究
3. 提案手法
4. データセット
5. 実験設定
6. 実験結果と分析
- 7. まとめ**

まとめ

◆まとめ

■提案

1. 「2 特徴ベクトル」の文書対応手法：Mean Poolingベクトルだけではなく、文書の系列ベクトルの中の第1文ベクトルも使用
2. 分割方式：重複あり固定長分割 (OFLS)と言語間の文書長(トークン数)の比率 ρ 付き重複あり固定長分割 (LD-OFLS)

■実験結果

- 重複あり「固定長分割」は「文単位分割」の性能を上回る
- 「固定長分割」は「文単位分割」より計算速度が速い
- 重複あり「固定長分割」でSentence Embeddingに基づく文書対応手法は、Baseline (機械翻訳によるDocalign) の性能を上回る

◆今後の展望

- 既存の公開データセットを用いた評価