

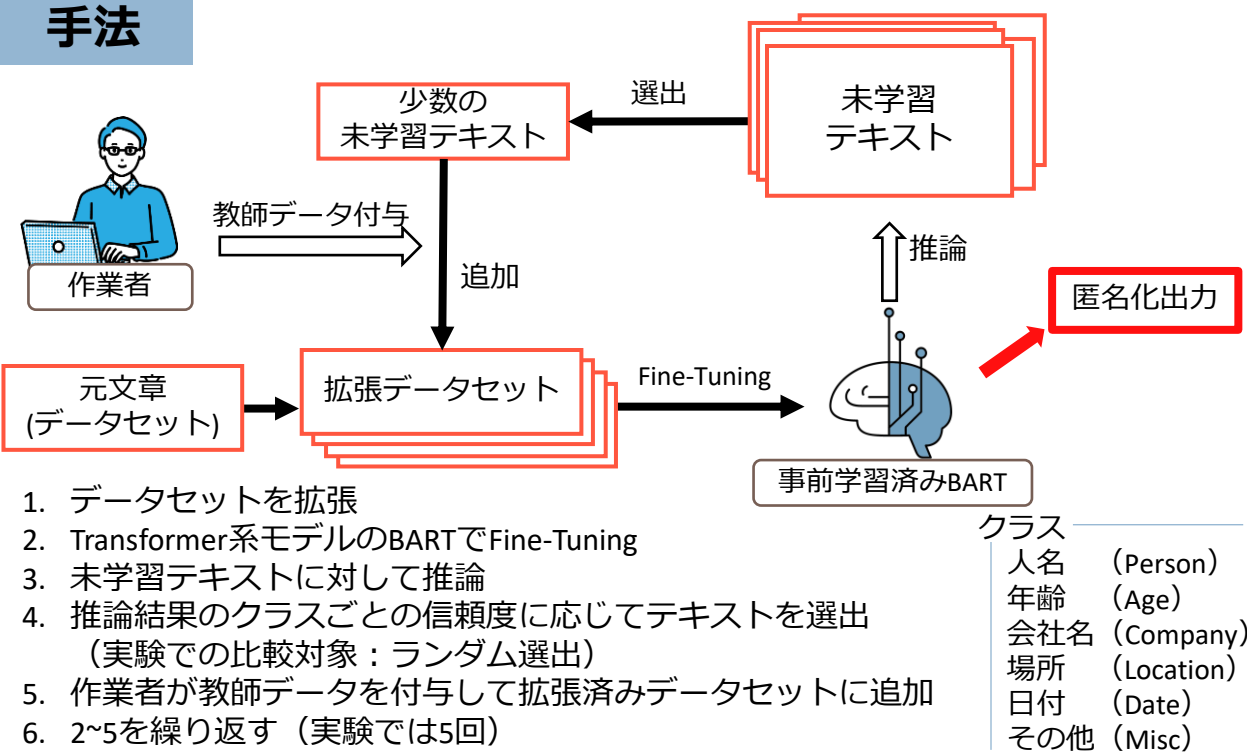
## 背景・目的

個人情報を保護するための匿名化処理がボトルネックとなり、判例や法律関係文書の一般公開が進んでいない  
単語の種類（クラス）ごとに登場数や識別難度に違いがある



- 機械学習モデルによる自動匿名化
- アクティブラーニングの導入による効率的な学習データの追加
- クラスごとの登場数・難度を是正するための学習

## 手法



- データセットを拡張
- Transformer系モデルのBARTでFine-Tuning
- 未学習テキストに対して推論
- 推論結果のクラスごとの信頼度に応じてテキストを選出 (実験での比較対象：ランダム選出)
- 作業者が教師データを付与して拡張済みデータセットに追加
- 2~5を繰り返す (実験では5回)

謝辞：本研究では、国立情報学研究所のIDRデータセット提供サービスにより弁護士ドットコム株式会社から提供を受けた「弁護士ドットコムデータセット」を利用した

## データセット

弁護士ドットコムデータセットを使用  
一部に仮想の個人情報を付与

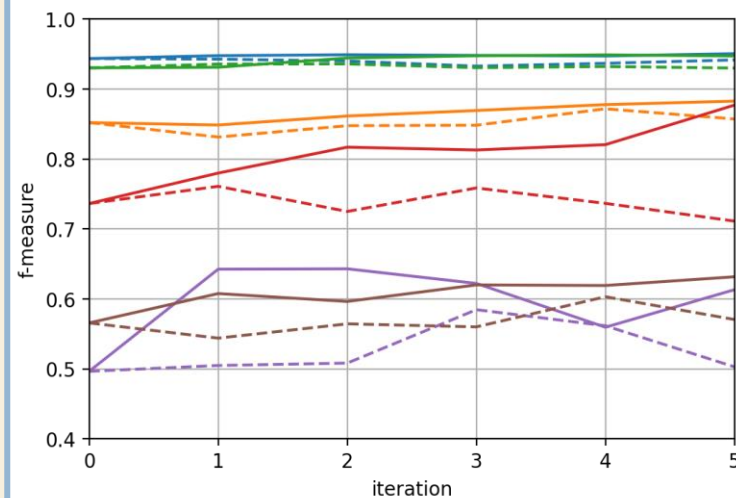
> もし18歳未満だったら逮捕や警察から呼び出しの可能性は高いですか？何度も22歳だと言われたのであれば、それを  
> 影山彰弁護士へ主張書を弁護士の方と相談しながら作成したいのですが、兵庫県の弁護士の方を紹介願えないでしょうか。  
野上様がクインライナー社に対し、「同僚に奥様がいて、お店を経営しているから保証人になります」という動機を表示し

弁護士への相談&回答テキスト、一部人名など個人情報を含む

## 結果

### 匿名化結果(強調表示)

> もし18歳未満だったら逮捕や警察から呼び出しの可能性は高いですか？何度も **b歳 Age** だと言われたのであれば、  
> **J氏 Person** 弁護士へ主張書を弁護士の方と相談しながら作成したいのですが、**β Location** の弁護士の方を紹介願え  
**N氏 Person** 様が **b社 Company** に対し、「同僚に **●○ Misc** がいて、お店を経営しているから保証人になります」



← アクティブラーニングによる学習の結果

人名・会社名で  
9割を超える精度

全クラスで精度  
向上、かつ  
ランダム追加  
より高い効果を  
確認