

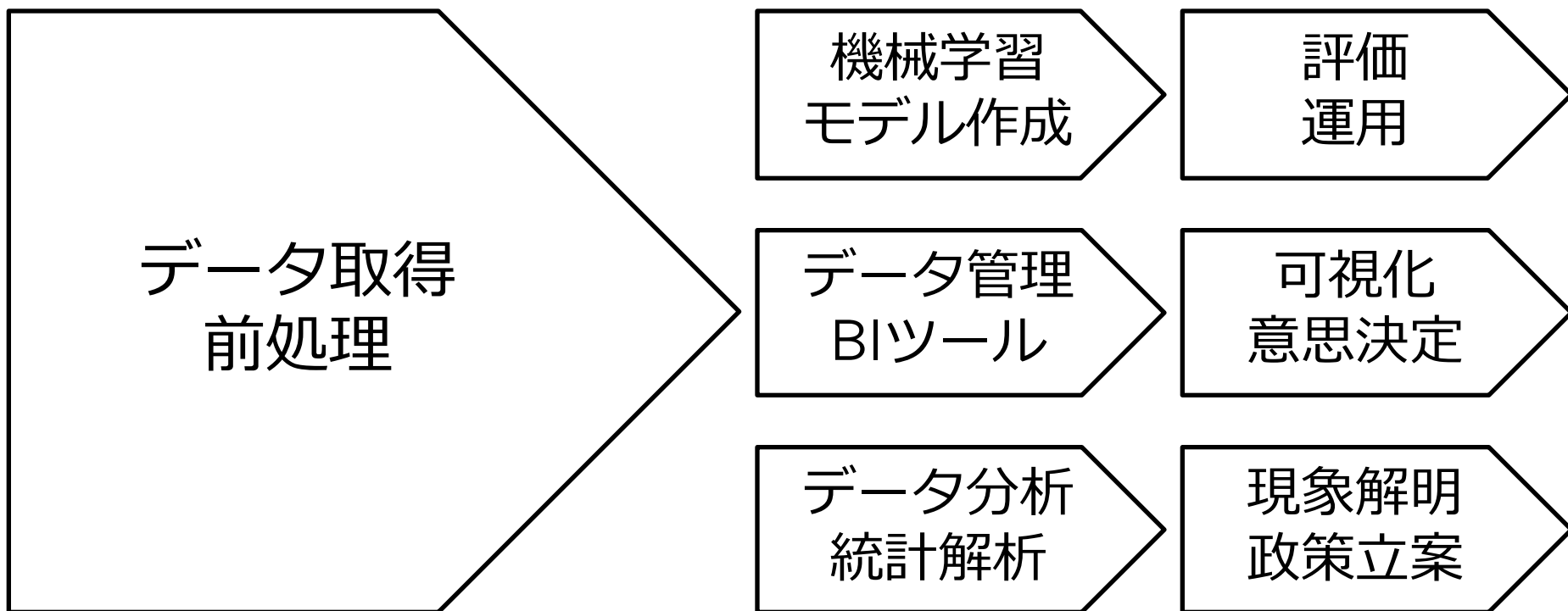
# BigGorillaの活動紹介

～データ前処理ノウハウの共有化に向けて～



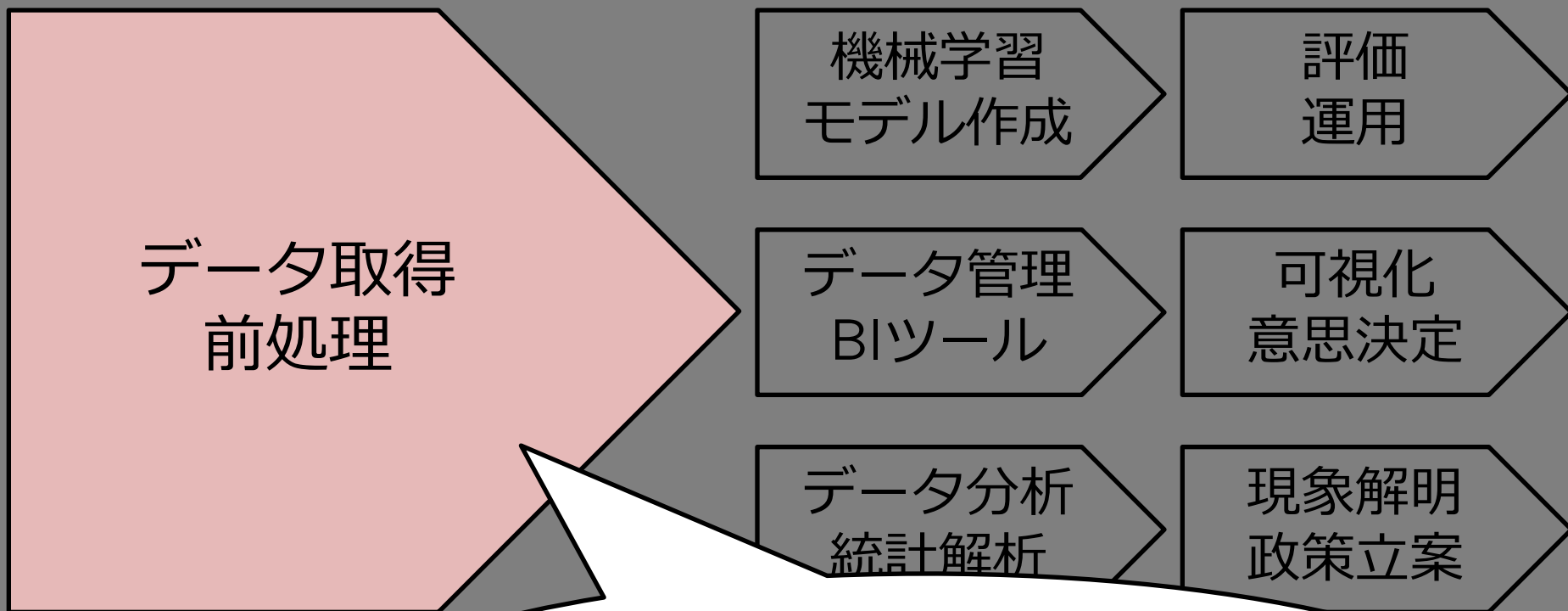
RECRUIT Institute of Technology 山下 宙元

# データ活用のフロー



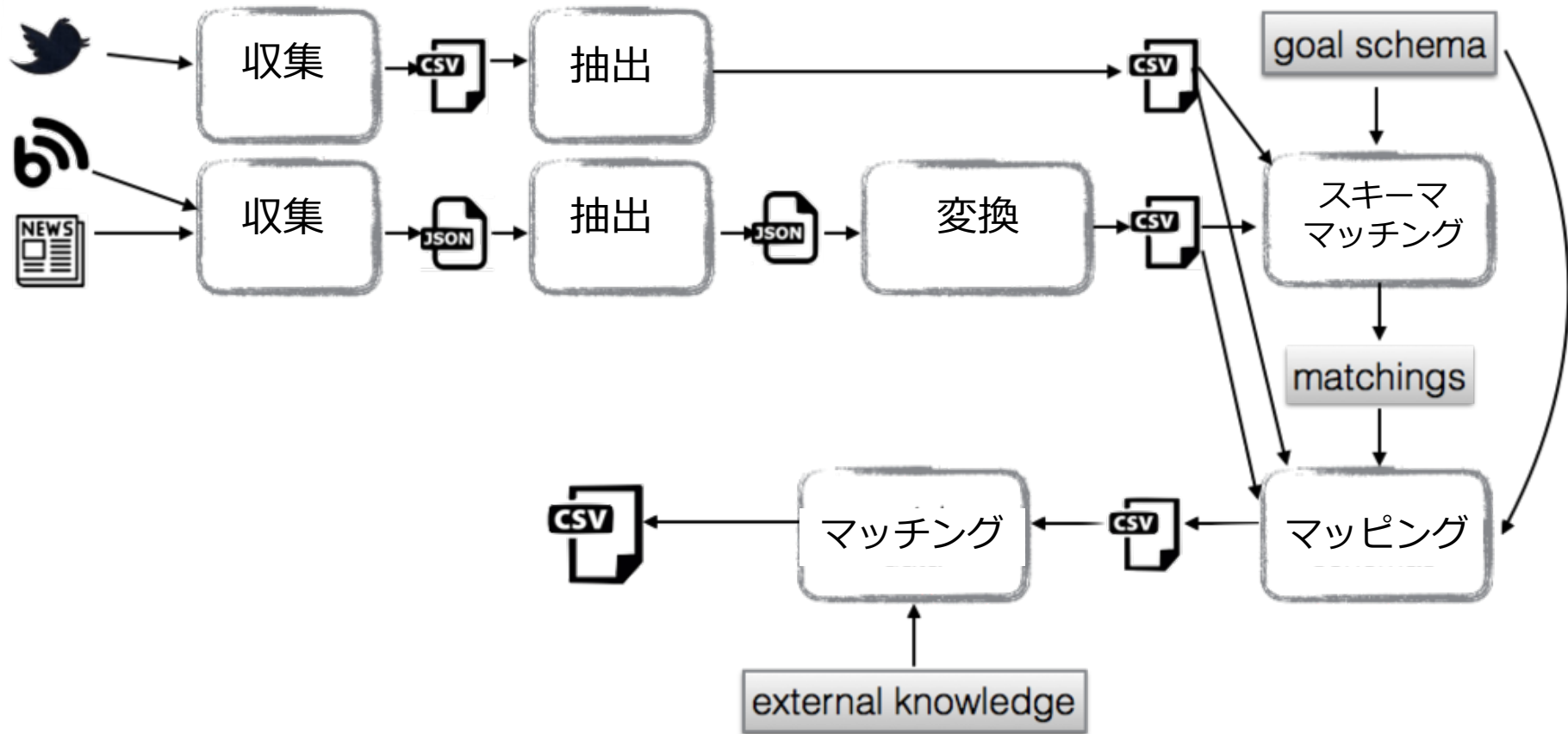
データ量の急速な増加、オープン分析ツールの増加により  
データの利活用が高まっているが…

# データ活用のフロー

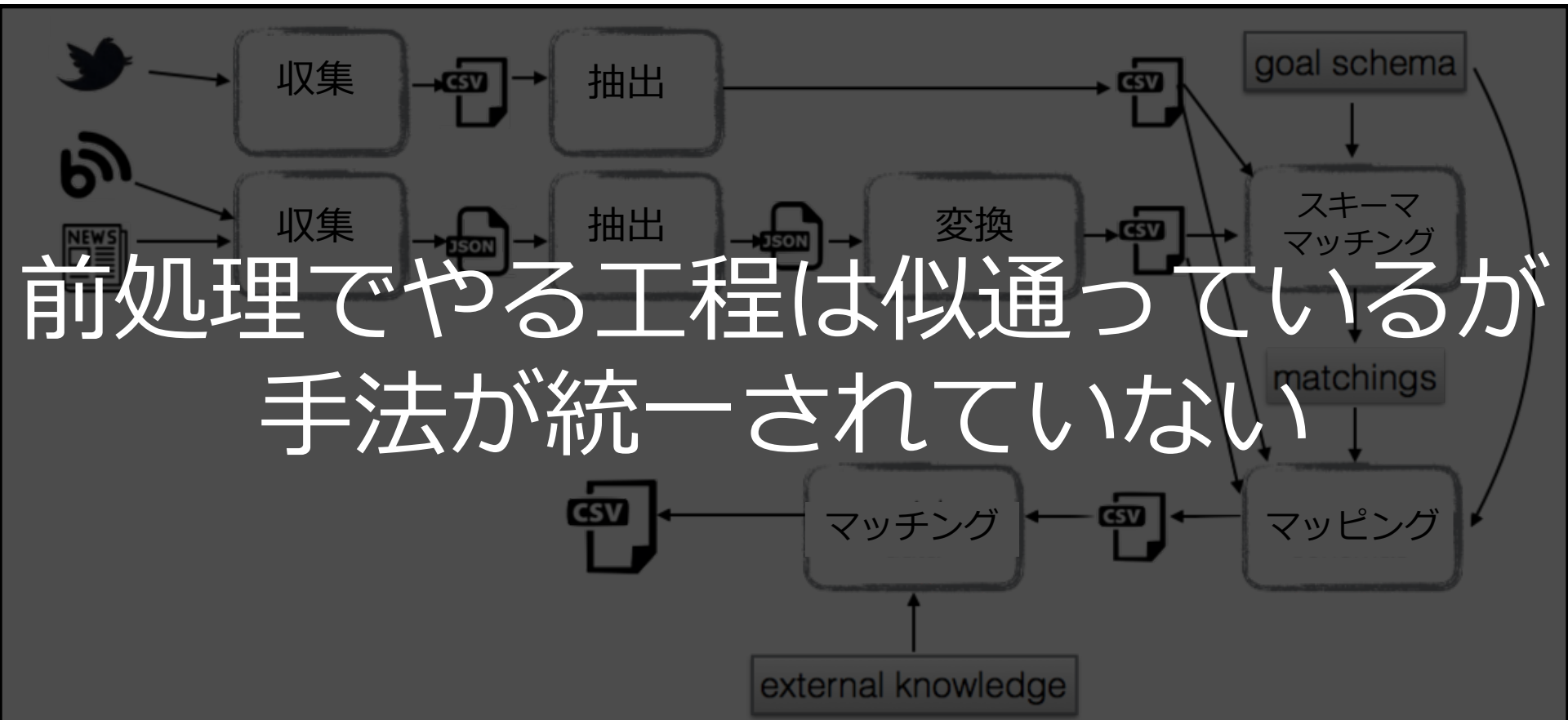


データ利活用において  
工数の8割は前処理

# 前処理には複数のタスクがある



# 前処理には複数のタスクがある



# BigGorillaのミッション



## ■ BigGorilla とは？

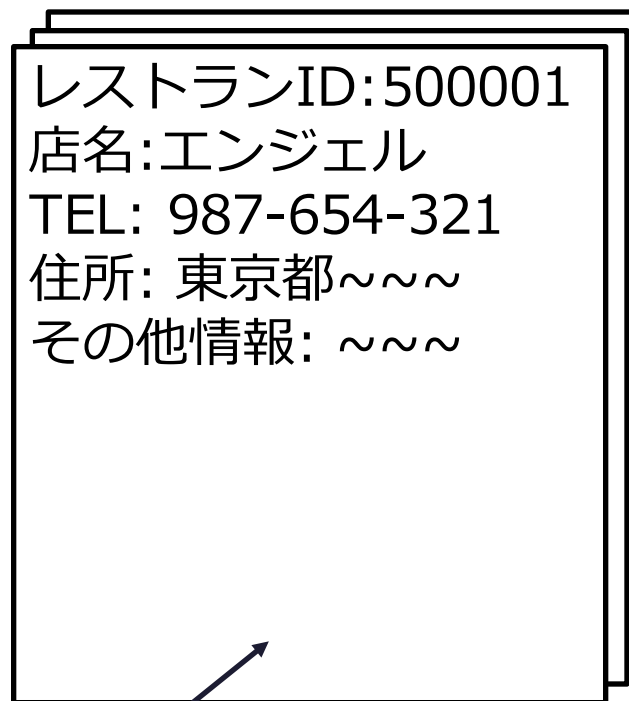
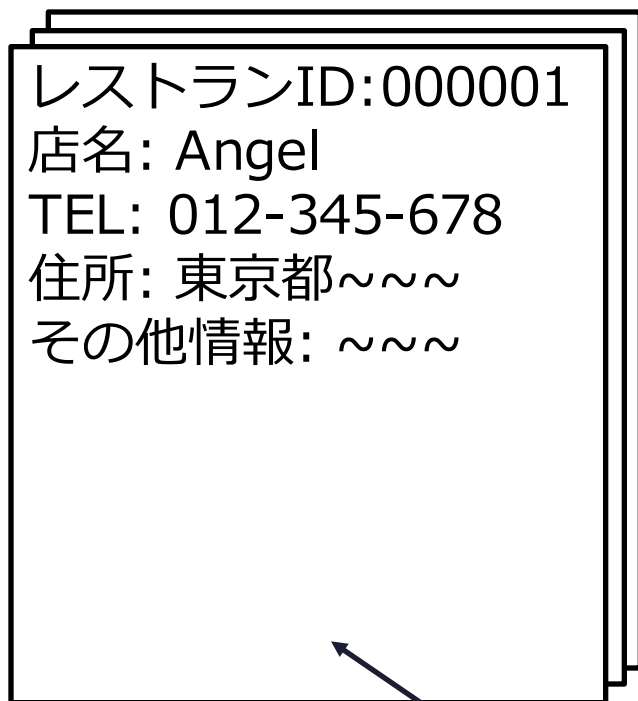
データ準備とデータ統合のためのオープンソース・エコシステムを開発/啓蒙する**オープンコミュニティ** (RITとWisconsin大学によりプロジェクト開始)

## ■ 活動の目的

データに向き合う人が直面する困難を解決するための**知識共有の場** (ソフトウェア/ワークフロー/データセット) となること

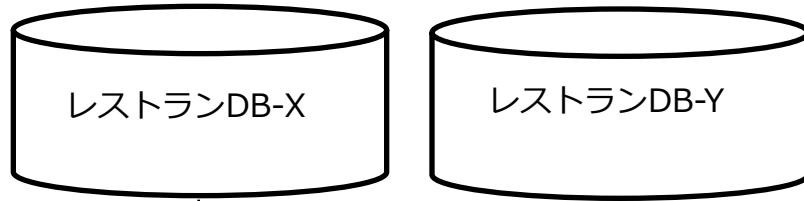
# ケーススタディ

～表記揺れのあるデータセットの名寄せ～

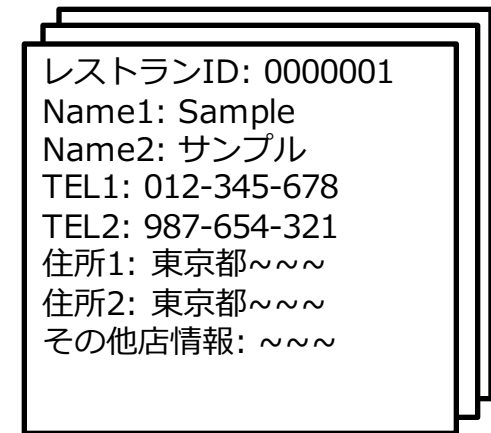
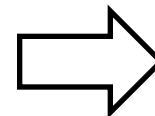
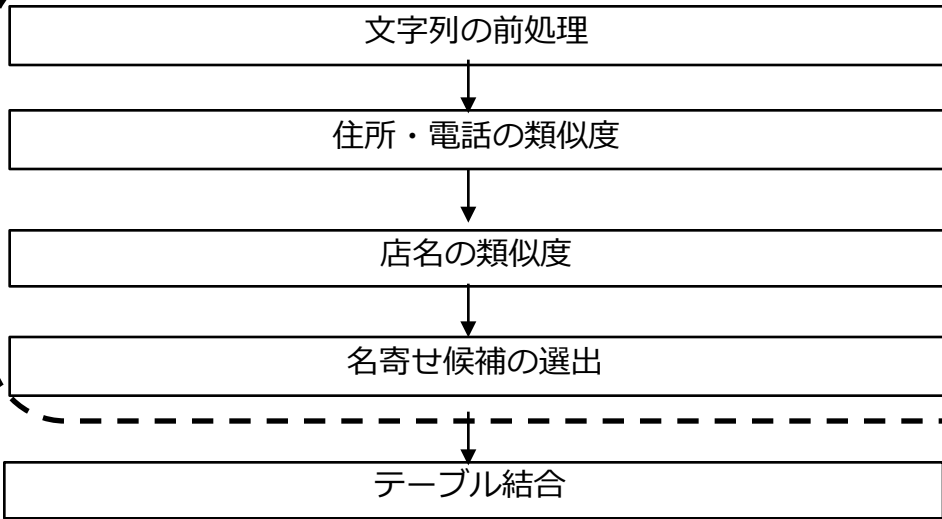


名寄せしたいが…  
人手で辞書をつくるのは面倒…

# 名寄せのアーキテクチャ



名寄せタスク





# 名寄せのアーキテクチャ

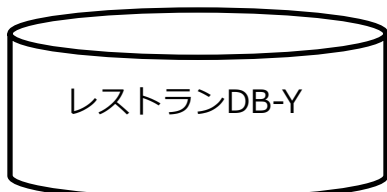
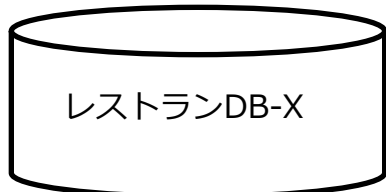
タスクの自動化

文字列の正規化

文字列の類似度

マッチング

データの結合



レストランID:000001  
店名: Angel  
TEL: 012-345-678  
住所: 東京都~~~  
その他情報: ~~~

レストランID:000001  
店名: エンジェル  
TEL: 987-654-321  
住所: 東京都~~~  
その他情報: ~~~

名寄せタスク

文字列の前処理

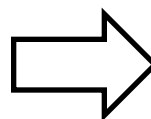
住所・電話の類似度

店名の類似度

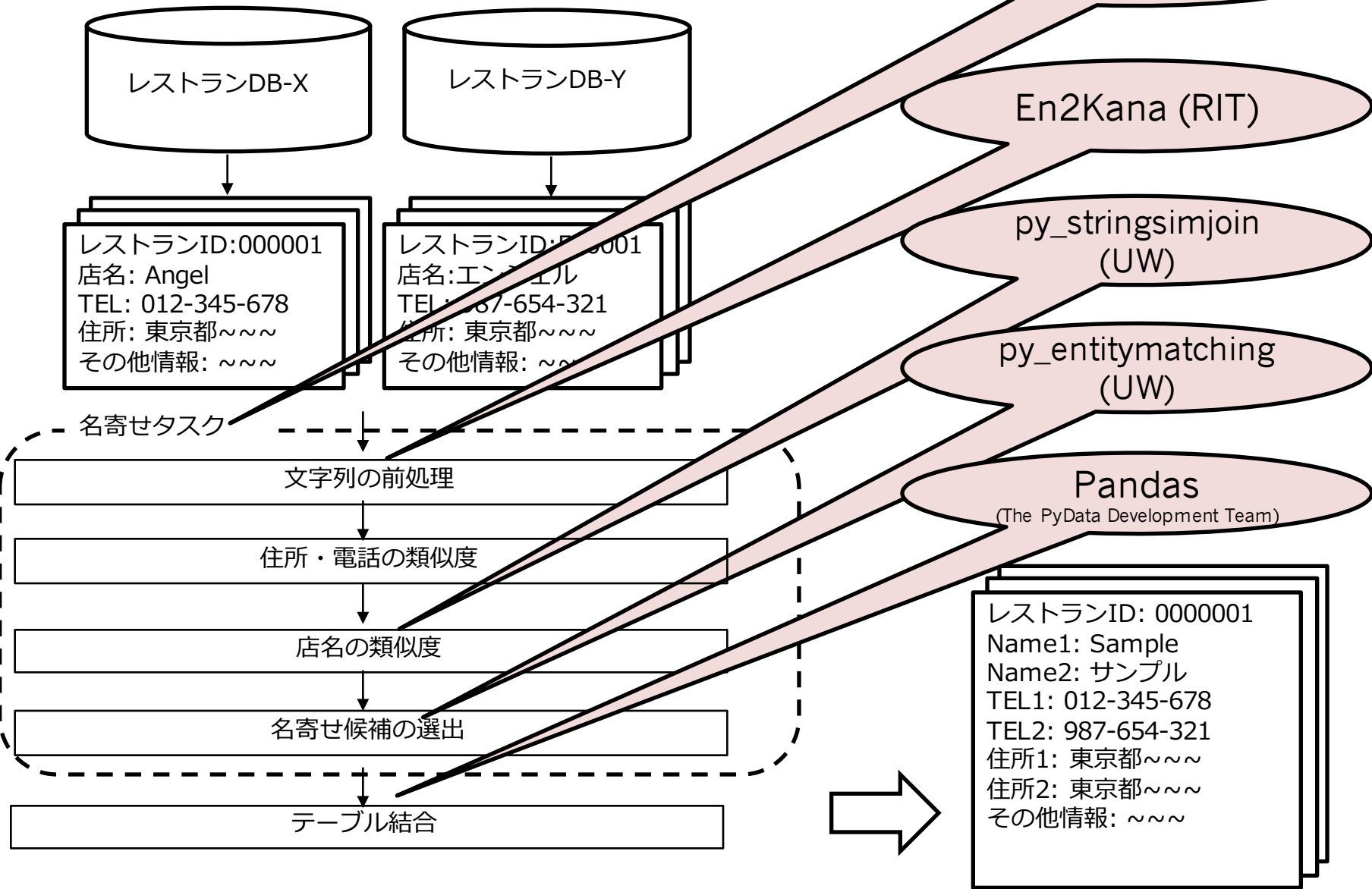
名寄せ候補の選出

テーブル結合

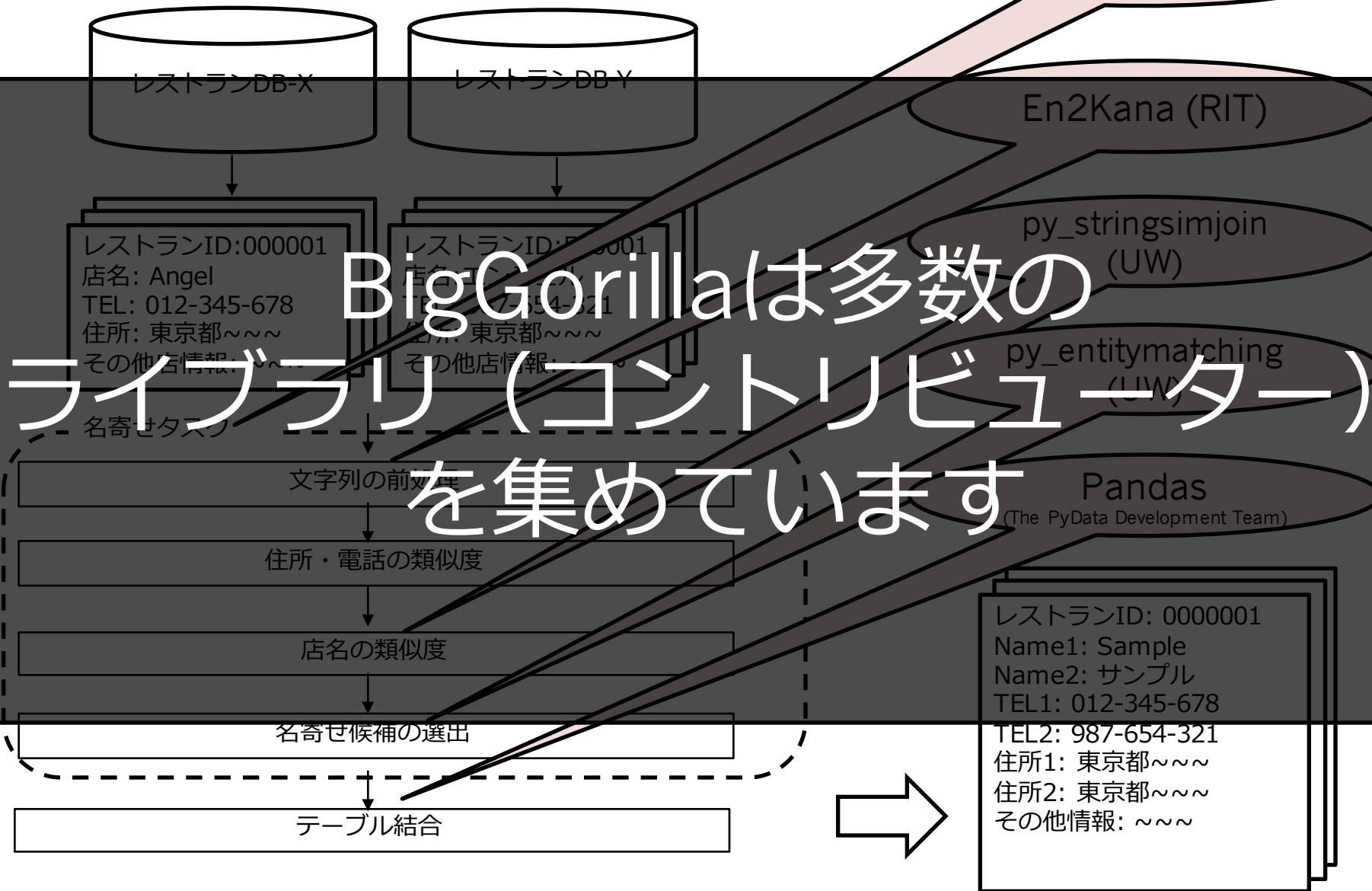
レストランID: 0000001  
Name1: Sample  
Name2: サンプル  
TEL1: 012-345-678  
TEL2: 987-654-321  
住所1: 東京都~~~  
住所2: 東京都~~~  
その他情報: ~~~



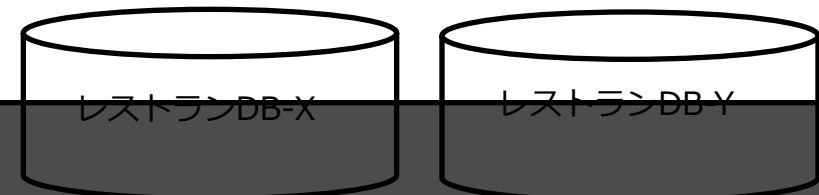
# 名寄せのアーキテクチャ



# 名寄せのアーキテクチャ



# 名寄せのアーキテクチャ



誰でもBigGorillaへ

ライブラリを投稿することができます

名寄せタスク

文字列の前処理

住所・電話の類似度

店名の類似度

名寄せ候補の選出

テーブル結合

Luigi (Spotify)

En2Kana (RIT)

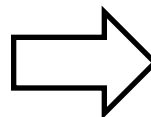
py\_stringsimjoin (UW)

py\_entitymatching (UW)

Pandas

(The PyData Development Team)

レストランID: 0000001  
Name1: Sample  
Name2: サンプル  
TEL1: 012-345-678  
TEL2: 987-654-321  
住所1: 東京都~~~  
住所2: 東京都~~~  
その他情報: ~~~



# En2Kana: 英語とカタカナの表記揺れを解決するモジュール

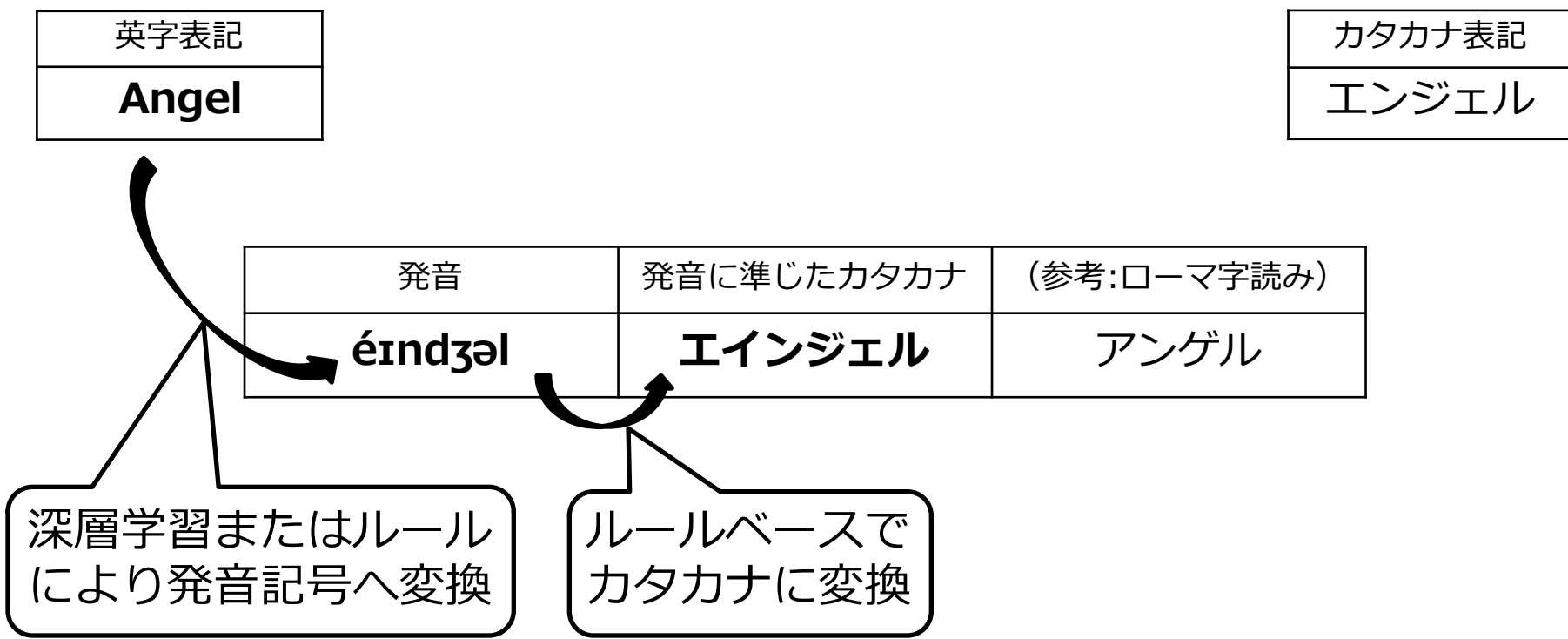
発音データを追加で用いることで  
文字情報のみでも名寄せが可能に



| 発音      | 発音に準じたカタカナ | (参考:ローマ字読み) |
|---------|------------|-------------|
| éɪndʒəl | エインジェル     | アングル        |

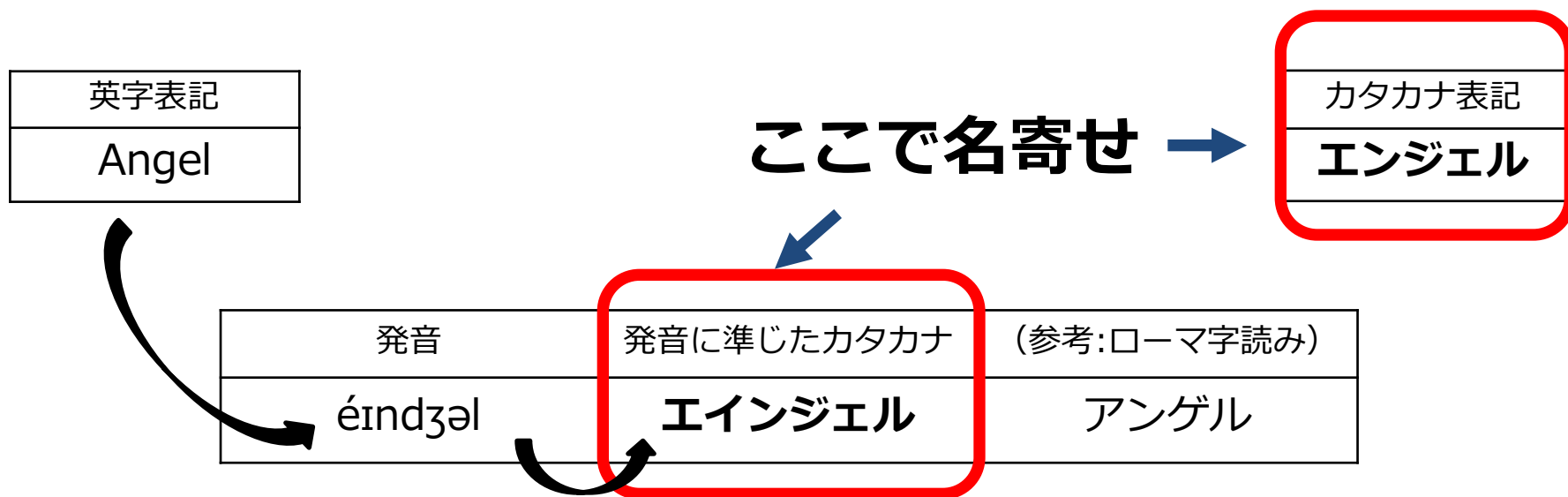
# En2Kana: 英語とカタカナの表記揺れを解決するモジュール

発音データを追加で用いることで  
文字情報のみでも名寄せが可能に



# En2Kana: 英語とカタカナの表記揺れを解決するモジュール

発音データを追加で用いることで  
文字情報のみでも名寄せが可能に



# テーブル名寄せの効果・影響

- 発音情報がないときよりも**10%以上の精度改善**
- 英語とカタカナの表記揺れの多い  
**10万件以上のデータテーブルを結合可能に**
- レストラン名だけでなく  
美容院名・会社名・人名・商品名など  
**他ドメインへも適用可能**
- **サーチエンジンのクエリ改善にも**



# BigGorilla その他のライブラリ

## BigGorilla:

データの前処理問題を解決するための  
**知識共有を目指したオープンコミュニティ**

## ライブラリ例:

Usagi: メタデータの検索、管理

KOKO: テキストデータから欲しい情報を抽出

他にも現在約**30**個のライブラリが集積  
(いつでも投稿お待ちしております)

# BigGorillaへのコントリビューション

## 1. ライブラリの紹介

既存のものでも自作でもOK

データの前処理に必要なものをご共有ください

## 2. 活用事例・ユースケースの紹介

どのような方法でどのような前処理問題を解決したか  
このような事例があったのでシェアしたいなど  
資料、ドキュメント、コード等をご共有ください

## 3. フォーラムへの投稿（意見・質問）

こういう前処理の問題で困っている  
こんな活用方法もあったなど

# BigGorillaへのコントリビューション

## 1. ライブラリの紹介

既存のものでも自作でもOK

データの前処理に必要なものをご共有ください

## 2. 活用事例（ソフトウェア）の問題を 一緒に解決しませんか？

どのようなライブラリを活用して  
どのような前処理問題を解決したか等の  
事例やコードをご共有ください

## 3. フォーラムへの投稿（意見・質問）

こういう前処理の問題で困っている  
こんな解決方法もあったなど