

クレジット:

UTokyo Online Education 学術俯瞰講義 2018 美馬秀樹

ライセンス:

利用者は、本講義資料を、教育的な目的に限ってページ単位で利用することができます。特に記載のない限り、本講義資料はページ単位でクリエイティブ・コモンズ 表示-非営利-改変禁止 ライセンスの下に提供されています。

<http://creativecommons.org/licenses/by-nc-nd/4.0/>

本講義資料内には、東京大学が第三者より許諾を得て利用している画像等や、各種ライセンスによって提供されている画像等が含まれています。個々の画像等を本講義資料から切り離して利用することはできません。個々の画像等の利用については、それぞれの権利者の定めるところに従ってください。



人工知能と自然言語処理を利用した 人文知の構造化

東京大学 大学総合教育研究センター
産業技術総合研究所人工知能研究センター
美馬 秀樹

講義内容

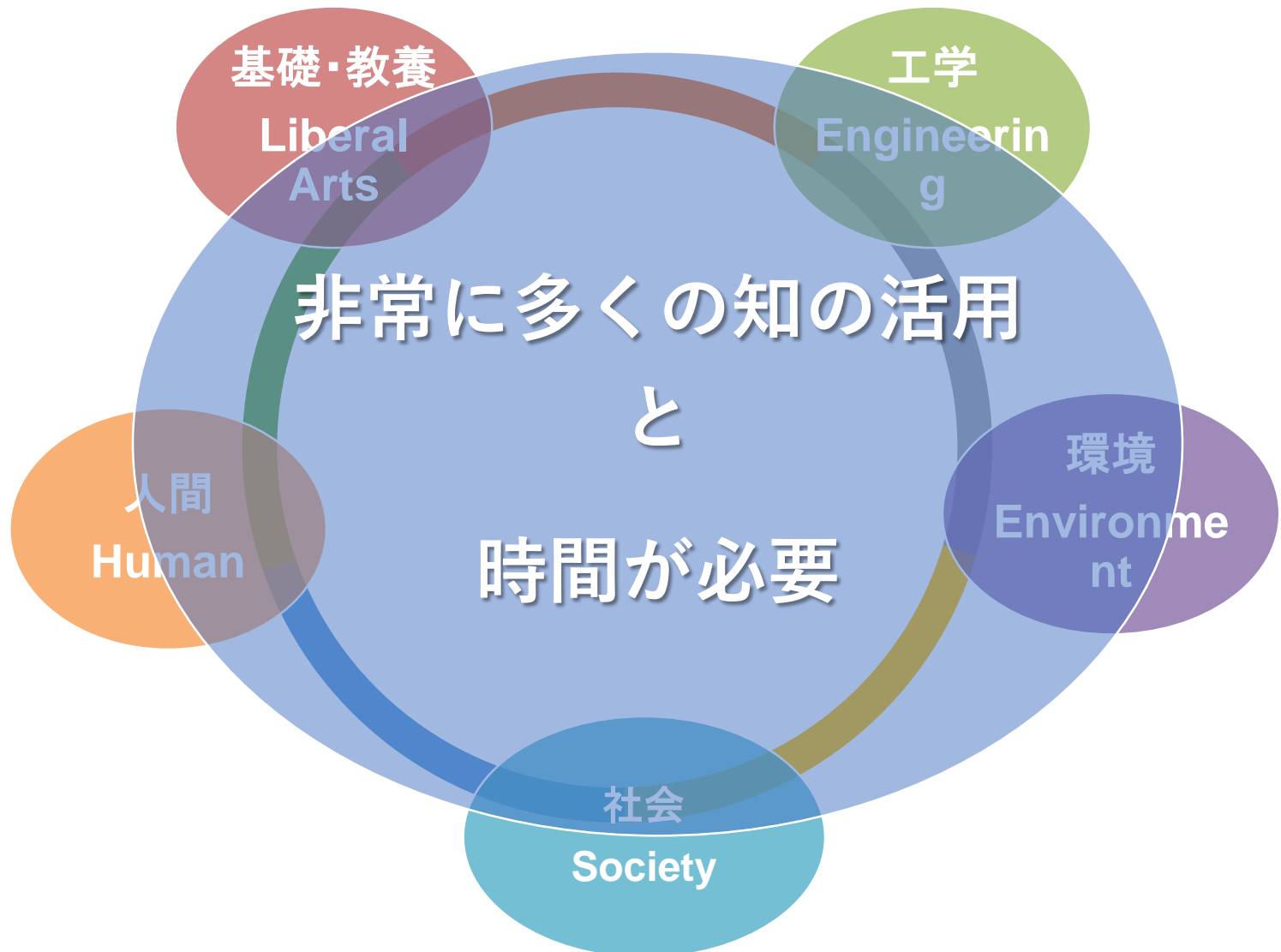
**I: 学際的教育、研究とデジタル
ヒューマニティーズ**

**II: 人工知能、自然言語処理、
デジタルアーカイブ化**

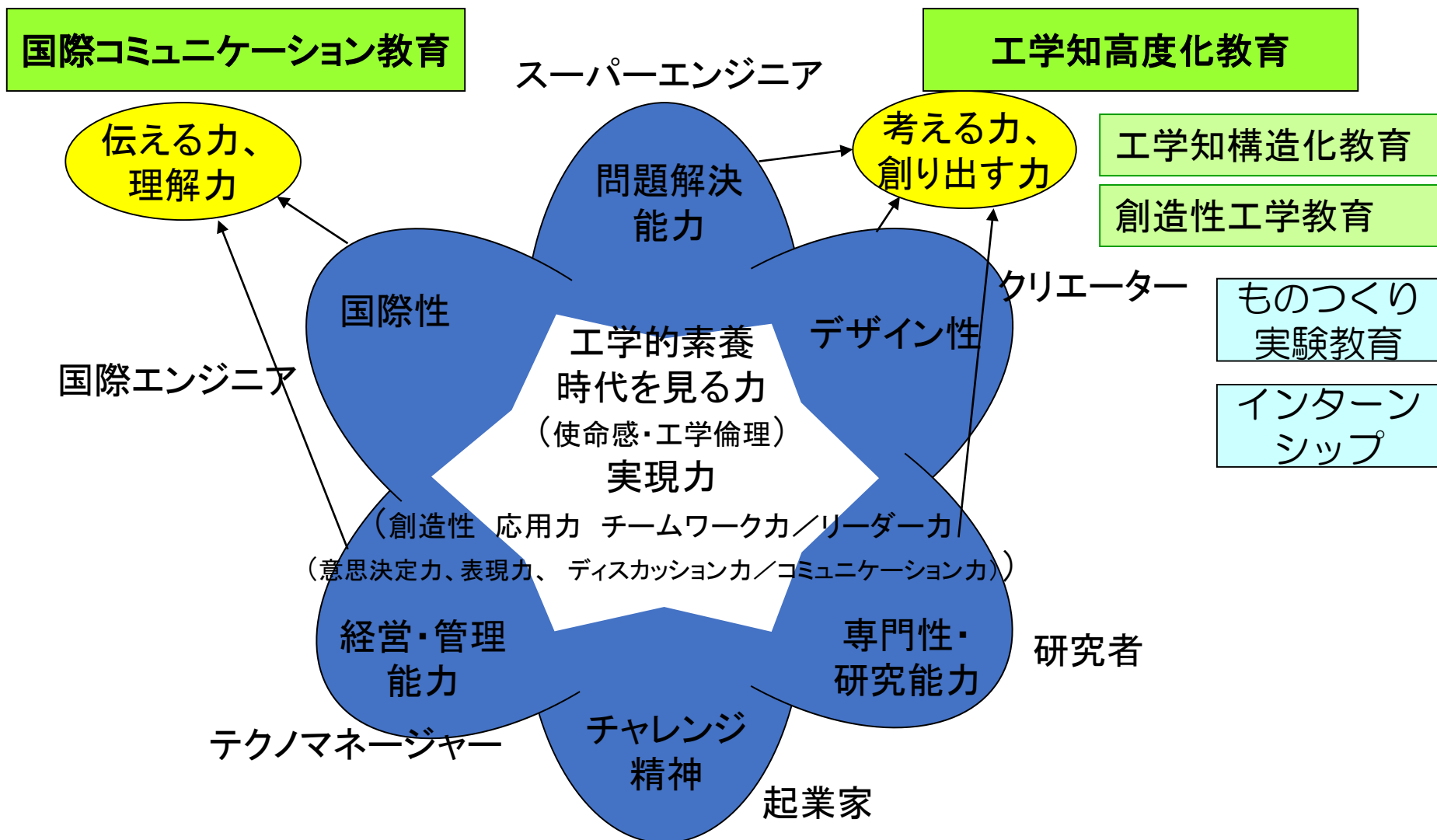
**III: 知の構造化の背景と岩波
書店「思想」のデジタル化**

IV: 現状と今後

ものづくりと学際性 Creative Engineering



多様でグローバルな人材のモデル



中島尚正編：「工学は何を目指すのか 東大工学部は考える」、東京大学出版会、2000 のP. 261に、加筆

多様でグローバルな人材のモデル

国際コミュニケーション教育

工学知高度化教育

スーパーエンジニア

伝える力、
理解力

問題解決
能力

考える力、
創り出す力

工学知構造化教育

創造性工学教育

クリエイター

ものづくり
教育

国際

社会の多様な課題に対して、
いかに多様な知識を活用するか？

経営・管理
能力

専門性・
研究能力

研究者

テクノマネージャー

チャレンジ
精神

起業家

学際的教育、研究の必要性

「分野、組織、時勢を越えて、
知を活用できるようにする」

文理融合、医工連携

FinTech、EdTech

etc.

背景

- 約1900万 / 毎月6万
 - 医学分野文献データベースに登録された文献数 (MEDLINE) / 毎月の増加数

学問における知識・情報の幾何級数的な増大
— 量の問題 —

- 約900 / 約12000
 - 東京大学 工学部の授業数 / 全学の授業数

学問の専門領域の細分化、従来の学問体系では対応できない複雑な問題の登場
— 質の問題 —

CVDの論文の数

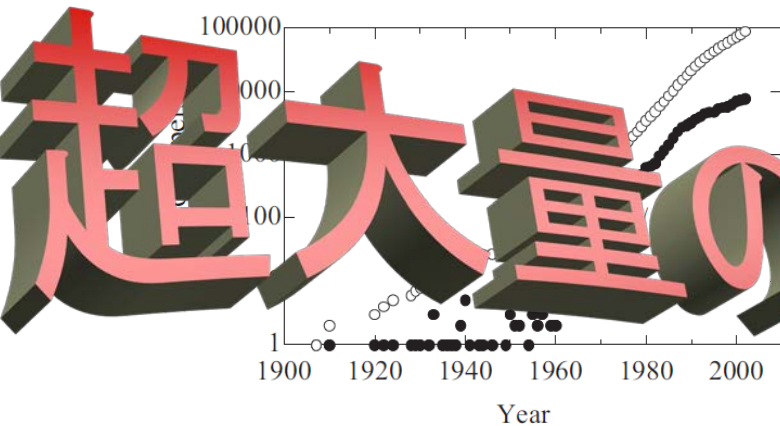
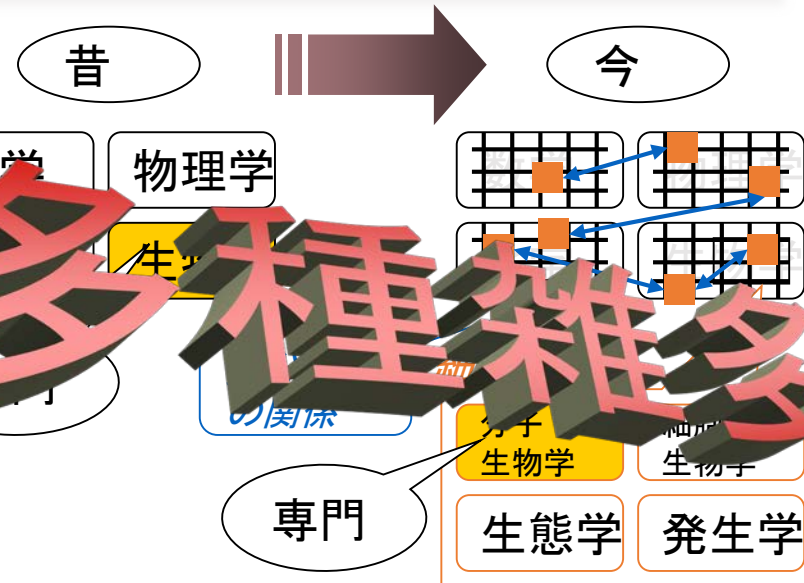


Fig. 1. Semi-log plot of the number of papers related to CVD. The closed circles are the number of papers published in that year. The open circles are the cumulative number of papers published up to and including that year.



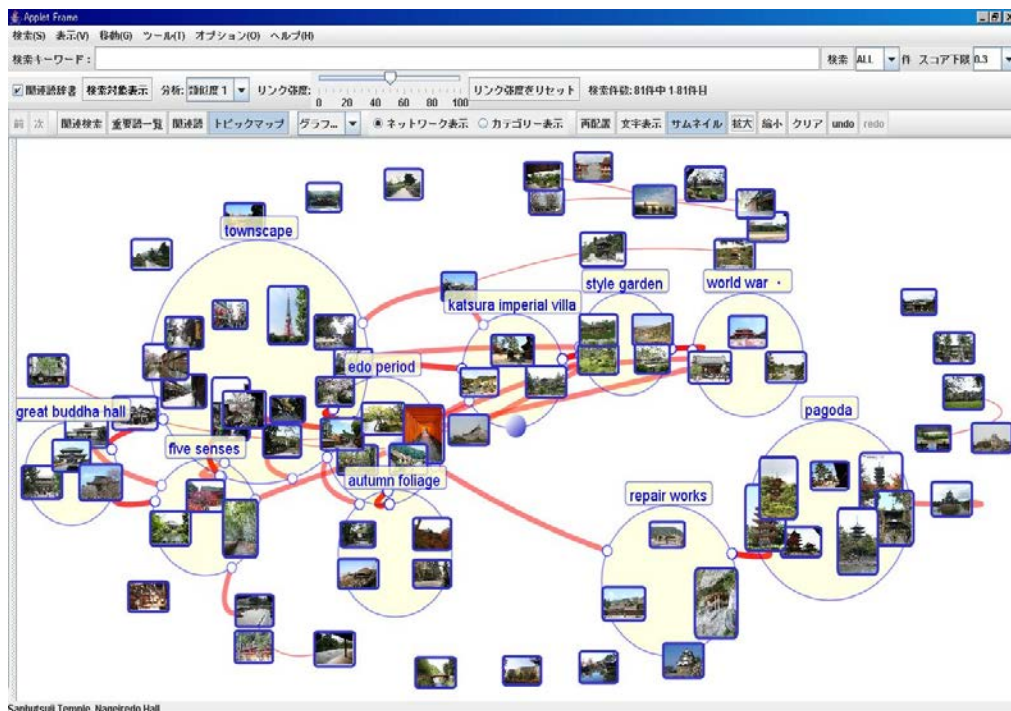
知の活用の技術

東京大学 美馬秀樹

1. 知識を収集・蓄積し、
2. 知識を分析・分類し、
3. 知識を可視化する。

テキストを解析し、
重要な概念を抽出し

関連性の認識を行う



文献1
ナノ粒子
チタニア
超臨界
サイズ
....

文献2
ナノ粒子
超臨界
発光
東北大
....

重要な概念の
重なりから 関
連性を推論

知の活用の技術

東京大学 美馬秀樹

1. 知識を収集・蓄積し、
2. 知識を分析・分類し、
3. 知識を可視化する。

情報技術の活用

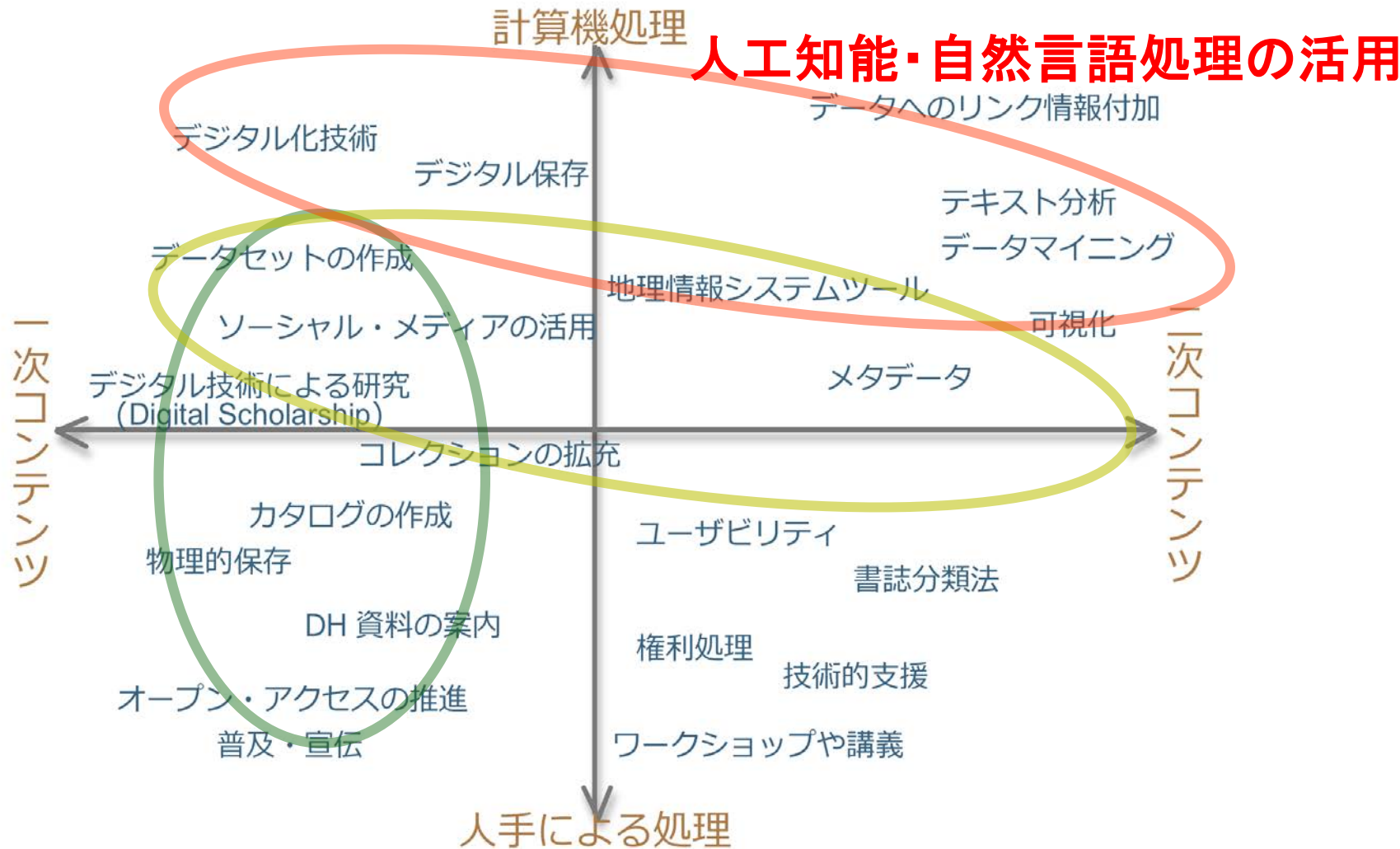


文献1
ナノ粒子
チタニア
超臨界
サイズ
....

重要な概念の
重なりから 関
連性を推論

文献2
ナノ粒子
超臨界
発光
東北大
....

DHの構造



Sula, Chris Alen, 2013, "Digital Humanities and Libraries: A Conceptual Model,"
Journal of Library Administration, 53(1) より

人工知能でできること

• 探索と最適化

- 路線検索
- 将棋、チェス
- 巡回セールスマン問題
- 最適配置

• 認識

- 五感の認識 + α
 - 音声認識
 - 文字認識、顔認識、物体認識
 - 味覚、におい、触覚
- 認知
 - 言葉の理解
 - 感情の理解

• 分類

- 行動モデリングと予測
- リコメンデーション

類推 →

似た情報を同様に処理

計算機の発展 →

大量の情報を高速に処理

最適化の例(本郷三丁目～駒場東大前)

- 経路探索
可能なルートから最適なものを選ぶ

- 評価軸
 - 運賃
 - 時間
 - 乗換回数
 - CO2量 etc.

本郷三丁目 ⇒ 駒場東大前 2007年4月10日 10時13分出発

表示順序を並び替える 所要時間 運賃 乗換回数

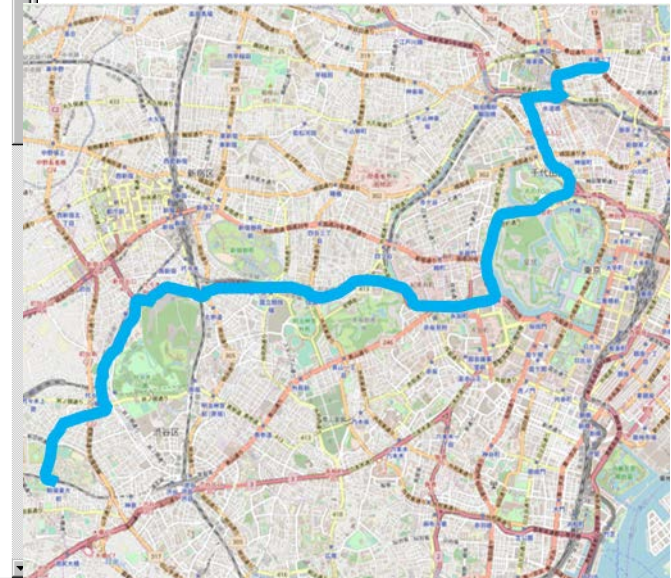
経路1 ● 10:17⇒10:52(35分) ¥310円 乗換回数:2回 CO2排出量:約226g(概算)

10:17発	本郷三丁目 [地図] [時刻表] [周辺検索]	190円
10:22着	東京メトロ丸ノ内線 荻窪行 2両目	05分
10:25発	大手町 [地図] [時刻表] [周辺検索]	
10:40着	東京メトロ半蔵門線 中央林間行 3・5両目	15分
10:49発	渋谷 [地図] [時刻表] [周辺検索]	120円
10:49着	京王井の頭線各停 吉祥寺行	03分
10:52着	駒場東大前 [地図] [時刻表] [周辺検索]	

経路2 ● 10:17⇒10:52(35分) ¥310円 乗換回数:2回 CO2排出量:約219g(概算)

10:17発	本郷三丁目 [地図] [時刻表] [周辺検索]	190円
10:32着	東京メトロ丸ノ内線 荻窪行 1・2・3・4・5・6両目	15分
10:35発	赤坂見附 [地図] [時刻表] [周辺検索]	
10:43着	東京メトロ銀座線 渋谷行 2・3両目	08分
10:49発	渋谷 [地図] [時刻表] [周辺検索]	120円
10:49着	京王井の頭線各停 吉祥寺行	03分
10:52着	駒場東大前 [地図] [時刻表] [周辺検索]	

経路3 ● 10:17⇒10:52(35分) ¥440円 乗換回数:3回 CO2排出量:約214g(概算)



(c)OpenStreetMap contributors

NAVITIME
<https://www.navitime.co.jp/>

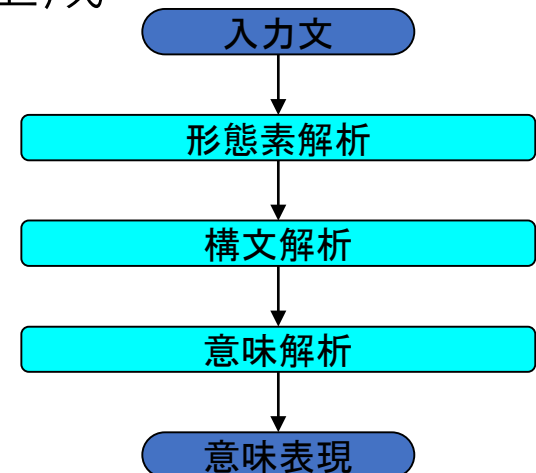
自然言語処理 (NLP)

• 計算機を用いて言語の理解を行う

- 形態素解析 – 単語（形態素）に区切る
- 構文解析 – 語構成、文の構成（主語、述語等）
- 意味解析 – 意味表現の生成

• アプリケーション

- 仮名漢字変換
- 機械翻訳
- 用語（概念）抽出
- 全文検索システム



計算機の発展 → 大量のテキストを高速に処理

自然言語処理と人工知能

- 「本郷三丁目から駒場東大前までの電車」

解析・構造化

現在地	行き先	手段
本郷三丁目	東大駒場前	電車

形態素解析

- 文を形態素（単語）に分け、品詞等の属性情報を同定する

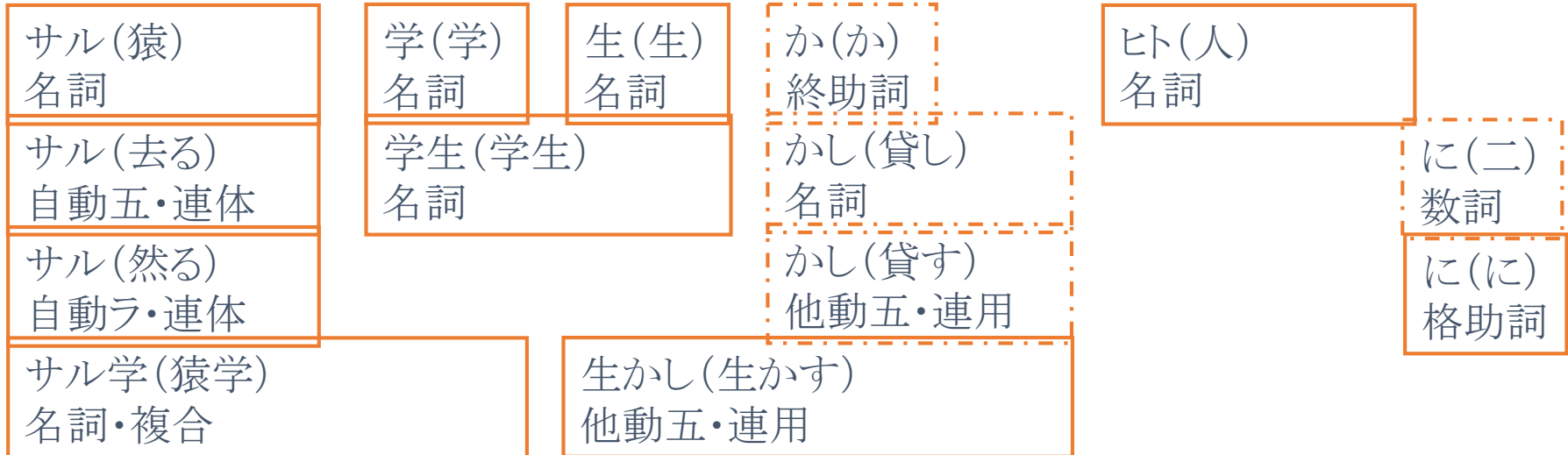
例：「本郷三丁目から駒場東大前までの電車」

形態素の表示

表記	よみ	品詞	基本形	全情報
本郷三丁目	ほんごうさんちょうめ	名詞	本郷三丁目	名詞,地名,*,本郷三丁目,ほんごうさんちょうめ,本郷三丁目
から	から	助詞	から	助詞,格助詞,*,から,から,から
駒場東大前	こまばとうだいまえ	名詞	駒場東大前	名詞,地名,*,駒場東大前,こまばとうだいまえ,駒場東大前
まで	まで	助詞	まで	助詞,副助詞,*,まで,まで,まで
の	の	助詞	の	助詞,助詞連体化,*,の,の,の
電車	でんしゃ	名詞	電車	名詞,名詞,*,電車,でんしゃ,電車

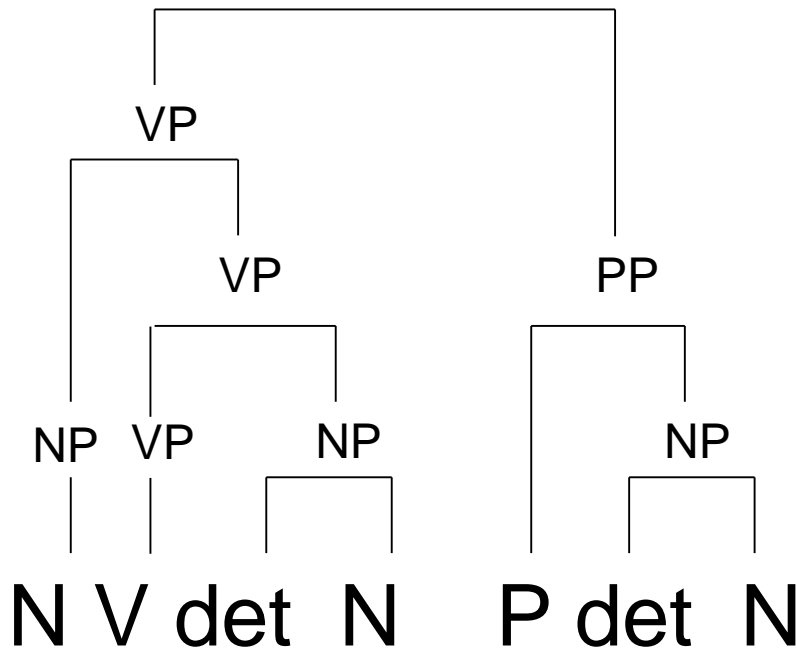
形態素解析の例

■ 例文(新聞の見出しから):
サル学生かしヒトに平和を



構文解析とあいまい性

- 文の統語的構造を同定する
 - 係り受け関係
- s • 例：「I saw a man in the park」



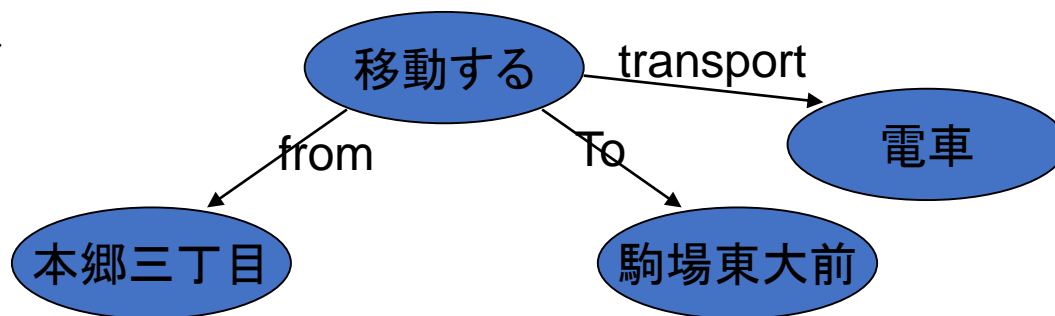
文法ルール

- S ← VP PP
- VP ← VP NP
- VP ← V
- NP ← det N
- NP ← N
- PP ← P NP

I saw a man in the park

意味解析

- 文の構造から意味表現を同定する（構造化）
- 文の意味とは
 - 意味ネットワーク



- フレーム

- (移動する
(from X?)
(to Y?)
(transport Z?)
)

実体化
→

(移動する-001
(from 本郷三丁目-001)
(to 駒場東大前-001)
(transport 電車-002)
)

構造化

現在地	行き先	交通手段
本郷三丁目-001	東大駒場前-001	電車-002

実検索

最適化の例(本郷三丁目～駒場東大前)

- 経路探索
可能なルートから最適なものを選ぶ

- 評価軸
 - 運賃
 - 時間
 - 乗換回数
 - CO2量 etc.

本郷三丁目 ⇒ 駒場東大前 2007年4月10日 10時13分出発

表示順序を並び替える 所要時間 運賃 乗換回数

経路1 10:17⇒10:52(35分) ¥310円 乗換回数:2回 CO2排出量:約226g(概算)

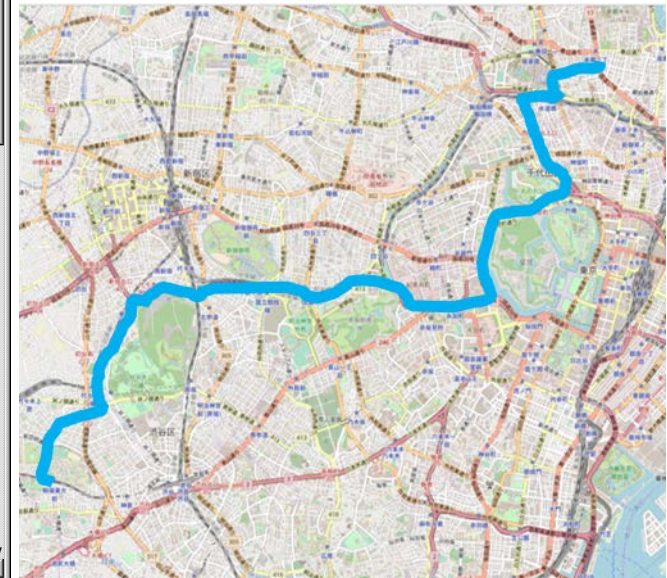
10:17発	本郷三丁目 [地図] [時刻表] [周辺検索]	190円
10:22着	東京メトロ丸ノ内線 荻窪行 2両目	05分
10:25発	大手町 [地図] [時刻表] [周辺検索]	
10:40着	東京メトロ半蔵門線 中央林間行 3・5両目	15分
10:49発	渋谷 [地図] [時刻表] [周辺検索]	120円
10:52着	京王井の頭線各停 吉祥寺行	03分
10:52着	駒場東大前 [地図] [時刻表] [周辺検索]	

経路2 10:17⇒10:52(35分) ¥310円 乗換回数:2回 CO2排出量:約219g(概算)

10:17発	本郷三丁目 [地図] [時刻表] [周辺検索]	190円
10:32着	東京メトロ丸ノ内線 荻窪行 1・2・3・4・5・6両目	15分
10:35発	赤坂見附 [地図] [時刻表] [周辺検索]	
10:43着	東京メトロ銀座線 渋谷行 2・3両目	08分
10:49発	渋谷 [地図] [時刻表] [周辺検索]	120円
10:52着	京王井の頭線各停 吉祥寺行	03分
10:52着	駒場東大前 [地図] [時刻表] [周辺検索]	

経路3 10:17⇒10:52(35分) ¥440円 乗換回数:3回 CO2排出量:約214g(概算)

NAVITIME
https://www.navitime.co.jp/



(c)OpenStreetMap contributors

「電車で駒場東大前まで本郷三丁目から行きたい」 → 同じ結果

意味のあいまい性

- 「apple」
 - 果物？ コンピュータ会社？
- 「Time flies like an arrow」
 - time – 時間、計る
 - fly – 飛ぶ、ハエ
 - like – 好き、～のように
- 「ぼくはキツネ」
 - 状況により様々な意味

文脈（コンテキスト）を理解 すること

- 「僕はキツネ！」



イラスト:いらすとや

岩波「思想」の構造化プロジェクト

- 構想

- 日本を代表する思想・哲学ジャーナル『思想』（1921年創刊）の80年分（約900号、約9,000論文、約16万ページ、かつ1986年までは活版印刷）
- **デジタル化と「知の構造化」**を行う

情報技術を活用して分析をしたい

デジタルのテキストデータがない

そもそもデジタルのデータがない

- 目的

- 『思想』という知の集積と20世紀日本の哲学・思想史を明らかにする
- 文献のデジタル化に関する方法論を確立する

K1B-33

5-19

思想

創刊號

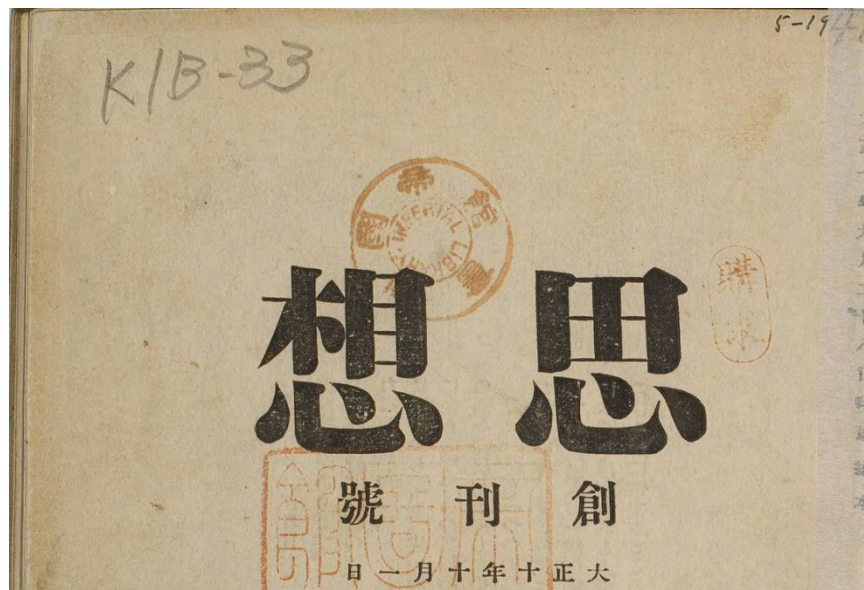
大正十年十一月一日

倉田百三	M. 和辻哲郎	M. 石原純	土居光知	桑木嚴翼	ケール	目次
父の心配……………(一)	世界見聞録……………(八)	原始基督教の文化的意義……………(五)	相對性原理の眞髓……………(四)	國民的文學と世界的文學……………(元)	流行の哲學思潮……………(三)	盛夏漫筆……………(一)

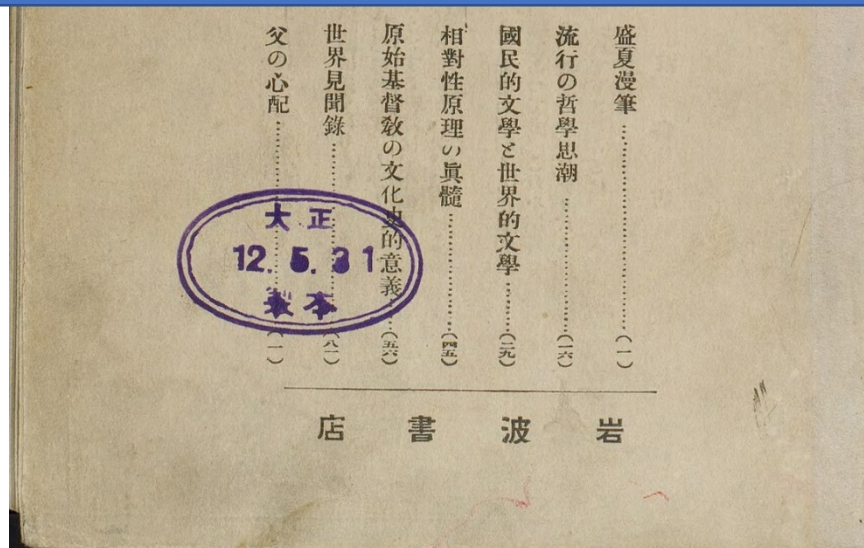
大正 12.5.31 本

岩波書店

『思想』創刊号、岩波書店、1921年



大量の古い文献を
どうやってデジタル化・テキスト化するか？



『思想』創刊号、岩波書店、1921年

テキストのデジタルアーカイブ

- コンピュータで処理できるテキストを構築
→ 検索や自然言語処理による解析が可能

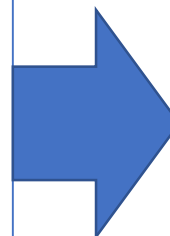


Photo by Ryan Hyde, from flickr
<https://flic.kr/p/89D449>
CC BY-SA 2.0

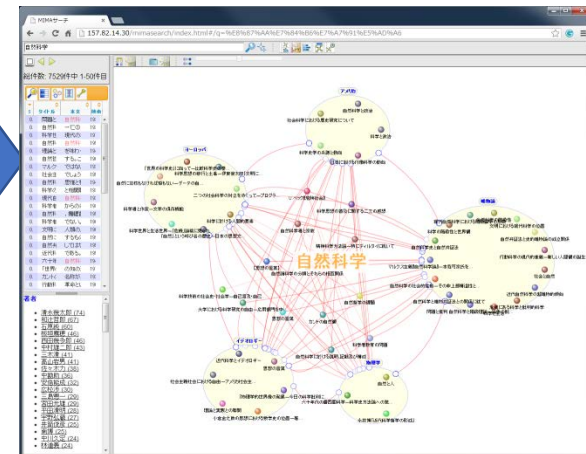


テキストデータ

するエリート待望論を持っていて、その点でニーチェに共鳴した面もあります。当時のアレヴィ・サークルには、二〇世紀前半ロマン・ロランとともにフランスの左翼知識人運動の中心を担ったアンリ・バルビュスや、のちの人民戦線の首班レオン・ブルム、そしてマルセル・ブルーストがいました。



分析システム



大規模なデジタルテキスト化

- **大量**の文献を**高精度**にテキスト化したい

大量の書籍・出版物

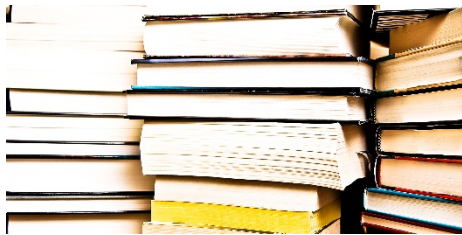


Photo by Ryan Hyde, from flickr
<https://flic.kr/p/89D449>
CC BY-SA 2.0



テキストデータ

するエリート待望論を持っていて、その点でニーチェに共鳴した面もあります。当時のアレヴィ・サークルには、二〇世紀前半ロマン・ロランとともにフランスの左翼知識人運動の中心を担ったアンリ・バルビュスや、のちの人民戦線の首班レオン・ブルム、そしてマルセル・ブルーストがいました。

大規模なデジタルテキスト化

大量の書籍・出版物

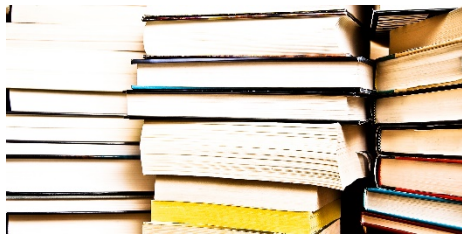


Photo by Ryan Hyde, from flickr
<https://flic.kr/p/89D449>
CC BY-SA 2.0



©いらすとや

人手入力

利点: ほぼ100%の精度

欠点: コストが高い

テキストデータ

するエリート待望論を持っていて、その点でニーチェに共鳴した面もあります。当時のアレヴィ・サークルには、二〇世紀前半ロマン・ロランとともにフランスの左翼知識人運動の中心を担ったアンリ・バルビュスや、のちの人民戦線の首班レオン・ブルム、そしてマルセル・ブルーストがいました。

大規模なデジタルテキスト化

大量の書籍・出版物

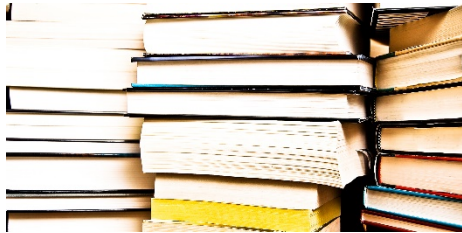


Photo by Ryan Hyde, from flickr
<https://flic.kr/p/89D449>
CC BY-SA 2.0



©いらすとや

音声認識

利点: 手入力よりも低コスト

欠点: 精度が低い

テキストデータ

するエリート待望論を持っていて、その点でニーチェに共鳴した面もあります。当時のアレヴィ・サークルには、二〇世紀前半ロマン・ロランとともにフランスの左翼知識人運動の中心を担ったアンリ・バルビュスや、のちの人民戦線の首班レオン・ブルム、そしてマルセル・ブルーストがいました。

大規模なデジタルテキスト化

大量の書籍・出版物

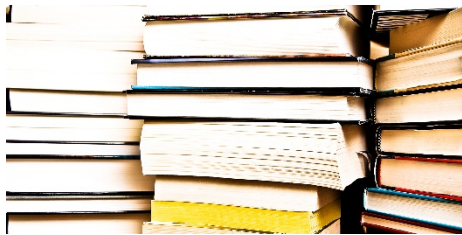
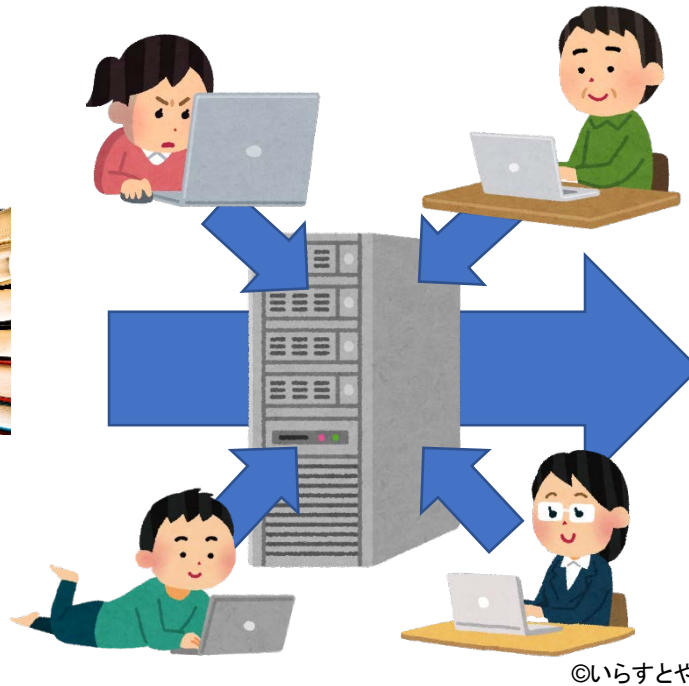


Photo by Ryan Hyde, from flickr
<https://flic.kr/p/89D449>
CC BY-SA 2.0



©いらすとや

手入力(クラウドソーシング)

利点: 高精度・低コスト

欠点: 質の統一が困難

データの管理

テキストデータ

するエリート待望論を持っていて、その点でニーチェに共鳴した面もあります。当時のアレヴィ・サークルには、二〇世紀前半ロマン・ロランとともにフランスの左翼知識人運動の中心を担ったアンリ・バルビュスや、のちの人民戦線の首班レオン・ブルム、そしてマルセル・ブルーストがいました。

クラウドソーシング例

NDLラボ 翻デジ

翻デジ2014
<http://lab.ndl.go.jp/dhii/omk2/>



アイテム	コレクションをブラウジング	Scripto	翻刻テキスト
------	---------------	---------	--------

はじめに

2014年3月24日

一般財団法人人文情報学研究所首席研究員・東京大学大学院情報学環特任准教授・
 立国会図書館研究員 [永崎研宣](#)

以下の説明の簡単な要約:

現状では、近代デジタルライブラリー(以下、近デジ)をテキストデータ化して色々便利に使えるようにするためのプラットフォームとなっております。関係者それぞれに多忙なかでの片手間の運営となっておりますのでおらかな気持ちで接していただけますと幸いです。徐々に色々改良・発展していく予定です。それから、[Mediawiki API](#)を使える方は少し面白いかもしれません。使い方はこのページの下の方にあります。

- このシステムは[日本デジタル・ヒューマニティーズ学会\(JADH\)](#)の分科会であるSIG-Transcribe JPが提供しています。
- このシステムは国立国会図書館次世代システム開発研究室が運営するサーバ上で動作しています。
- このシステムはジョージ・メイソン大学で開発された[Omeka](#)([Omekaの説明の日本語訳](#))というGLAM向けコンテンツマネジメントシステム(CMS)を核として構築されています。データの蓄積はMediawikiに行なわれるようになっていきます。
- このシステムでは、各資料のページ毎にデジタル翻刻を行っていただけますと、デジタル翻刻テキストの各頁から近デジの当該頁に自動的にリンクが張られるようになっており、翻刻テキストの頁から容易に実際の本文を画像で確認できるようになっています。したがって、翻刻テキストとしての正確性がそれほど高くなくとも検索用途等としてある程度利用可能なものを作成することができるということになっております。
- このシステムは以下のことを**目標としています**。
 - 近代デジタルライブラリー等の画像でしか提供されていない日本語

最近追加されたアイテム

[道後温泉誌](#)

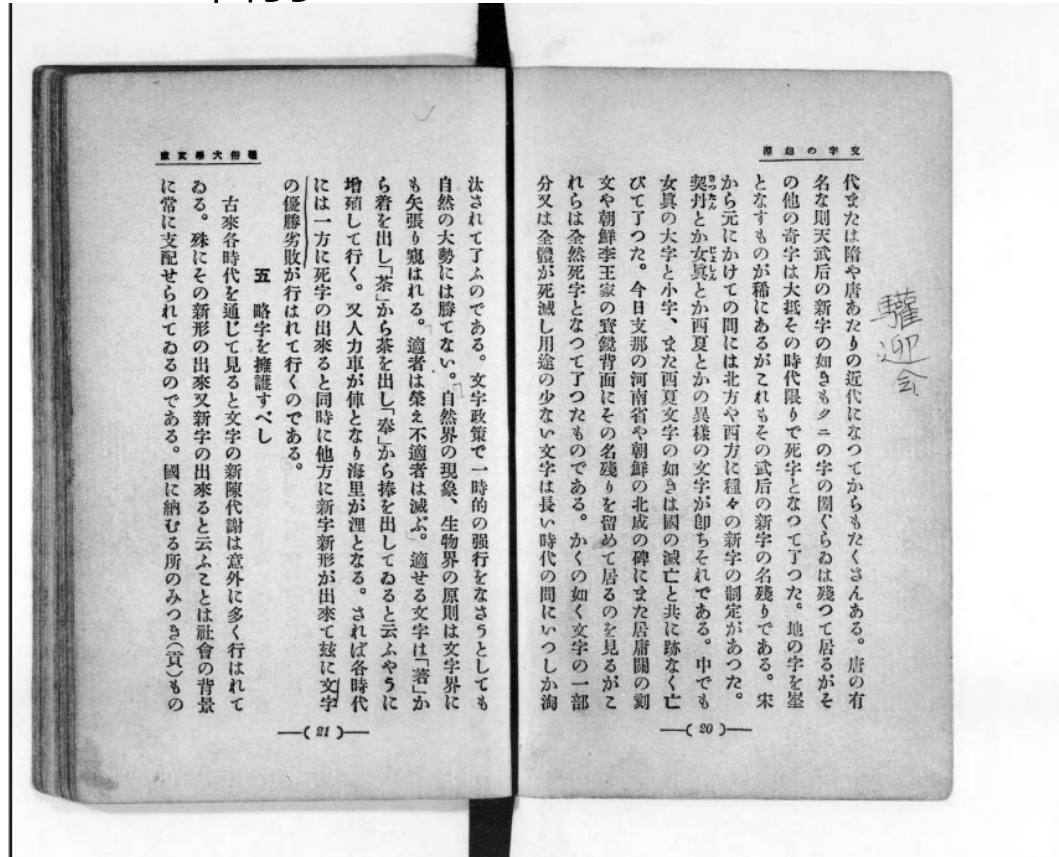
[天竺徳兵衛](#)

[あさどり](#)

[すべてのアイテムを見る](#)

クラウドソーシング例

NDLラボ 翻デジ



[« previous page](#) | [next page »](#) | [show discussion](#)

国立国会図書館デジタルコレクションより引用

You don't have permission to transcribe this page.

この頁を電子翻刻 [履歴]

新字旧字混在 ▾ 仮名遣い混在 ▾ Wikiのみ ▾

代または隋や唐あたりの近代になつてからもたくさんある。唐の有名な則天武后の新字の如きもクニの字の固くらは残つて居るがその他の奇字は大抵その時代限りで死字となつて了つた。地の字を鑿となすものが稀にあるがこれもその武后の新字の名残りである。宋から元にかけての間には北方や西方に種々の新字の制定があつた。契丹(きつた

大規模なデジタルテキスト化

大量の書籍・出版物

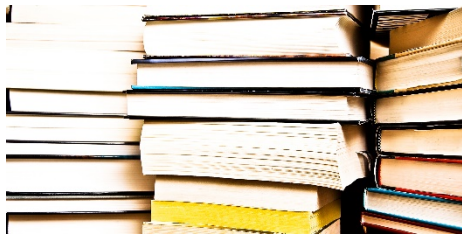


Photo by Ryan Hyde, from flickr
<https://flic.kr/p/89D449>
CC BY-SA 2.0



スキャン+ 文字認識

利点: コストが低い

短時間で処理可能

欠点: 精度が低い

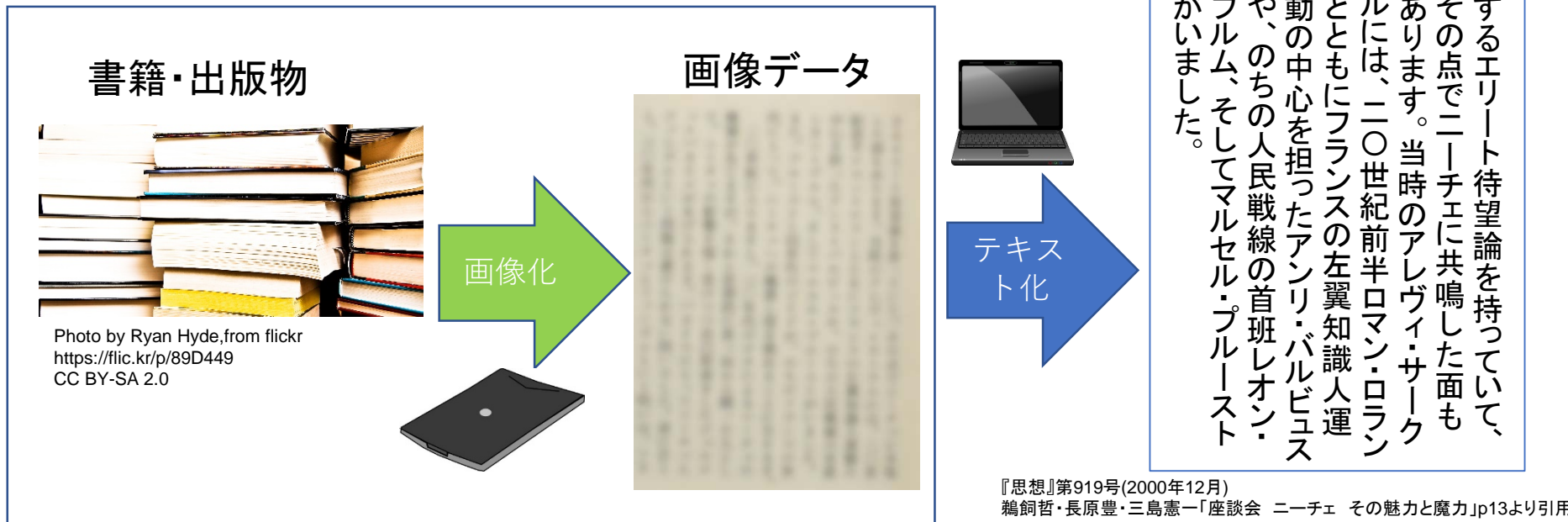
→技術発展により向上が望める

テキストデータ

するエリート待望論を持っていて、その点でニーチェに共鳴した面もあります。当時のアレヴィ・サークルには、二〇世紀前半ロマン・ロランとともにフランスの左翼知識人運動の中心を担ったアンリ・バルビュスや、のちの人民戦線の首班レオン・ブルム、そしてマルセル・ブルーストがいました。

デジタルテキスト化

- 印刷された文字をコンピュータで処理可能なデジタルテキストに変換
 - 画像化: スキャン
 - テキスト化: OCR(Optical Character Recognition)



デジタル画像化

- 対象物のスキャンによりデジタルデータ化
 - スキャナの種類(スキャン形態による分類)
 - フラットベッドスキャナ
 - 本をガラスの台に固定し、
下から光を当てて読み取る
 - シートフィードスキャナ
 - 自動原稿送り装置で
原稿を送り、読み取る
 - カメラスキャナ(スタンドスキャナ)
 - 原稿を専用の台に置き、カメラで撮影
 - 本の非破壊スキャンが可能

著作権等の都合により
ここに挿入されていた画像を
削除しました

フラットベッドスキャナの画像



<http://cweb.canon.jp/imageformula/lineup/p215ii/index.html>
Canon DR-p215

カメラスキャナ



自動的に本をめくり
上にある2台のカメラで撮影

本を痛めることなく
自動でスキャンが可能

キルタス社製BookScanner(APT BookScan 2400RA)

出典: 株式会社ブライトホープhttp://www.bright-hope.co.jp/cnt_02.html

カメラスキャナ？

- Scansnap SV600

- 実はカメラではなくラインセンサーでスキャン
- ラインセンサーとライトが首を振る
- 機能としてはカメラカメラスキャナと同様



FUJITSU
ScansnapSV600
充実したソフトウェア
<http://scansnap.fujitsu.com/jp/product/sv600/function.html>

Book Turner

CASIO

大切な本やノートを電子書籍化

電子書籍化支援システム

BOOK TURNER

ブックターナー



[裁断不要] [ページめくりをアシスト]

BT-100

購入前のご注意

- ・本製品は、完全自動で書籍を電子化するブックスキャナーではありません。使用者の補助が必要ですので、最初にトレーニングしてから、大切な本などを撮影することを推奨します。
- ・本やノートの紙質、厚さ、製本状態によっては、本製品のページめくり動作がうまくできない場合があります。その場合は、半自動モードもしくは手動モードで撮影してください。
- ・本製品を使用する場合は、著作権法を遵守してください。
- ・本製品を使用するには、タブレットまたはスマートフォンと専用アプリが必要です。
- ・本製品には、タブレットやスマートフォンは含まれません。

テキストデータの取得

- 印刷された文字をコンピュータで処理可能なデジタルテキストに変換
 - 画像化: スキャン
 - テキスト化: OCR(Optical Character Recognition)

『思想』第919号(2000年12月)
 鴉飼哲・長原豊・三島憲一「座談会 ニーチェ その魅力と魔力」p13より引用

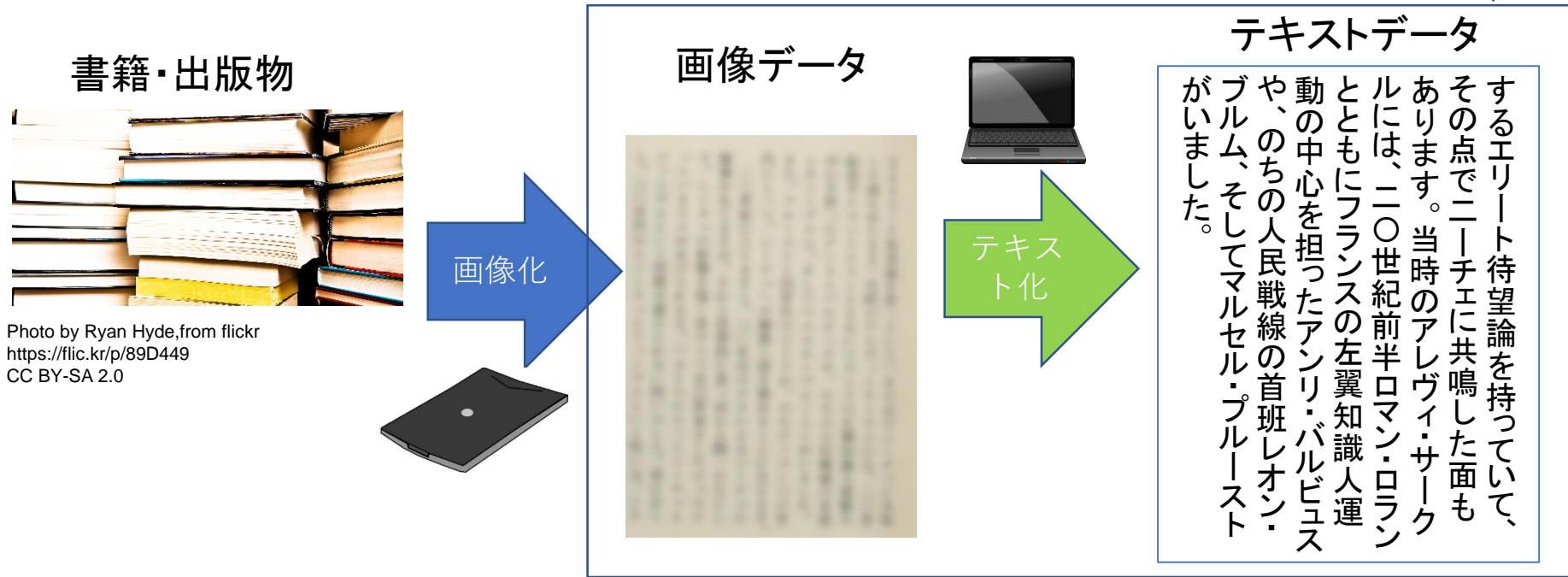
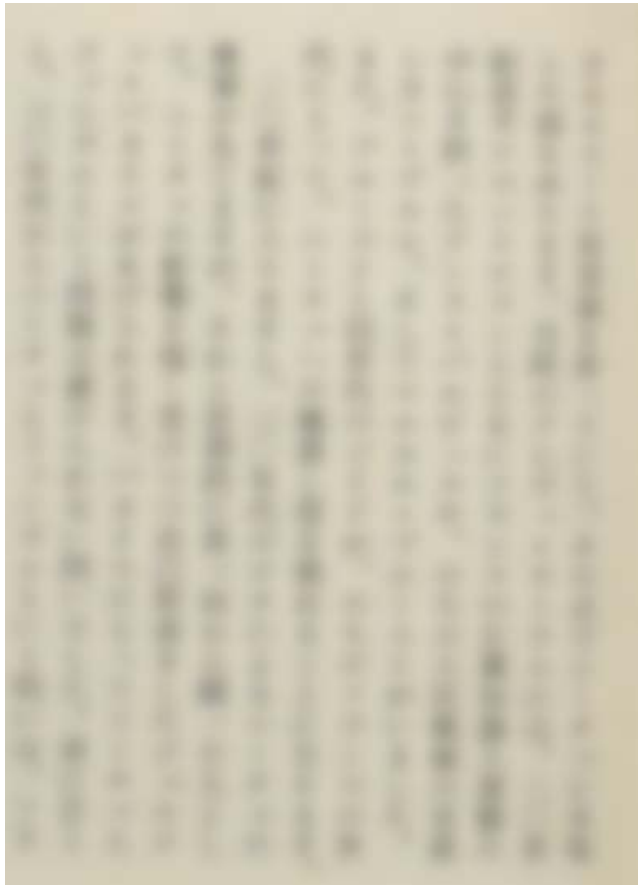


Photo by Ryan Hyde, from flickr
<https://flic.kr/p/89D449>
 CC BY-SA 2.0

デジタルテキスト化: OCR(Optical Character Recognition)

『思想』第919号(2000年12月)
鴉飼哲・長原豊・三島憲一「座談会 ニーチェ その魅力と魔力」p13より引用

- デジタル画像から文字を認識する



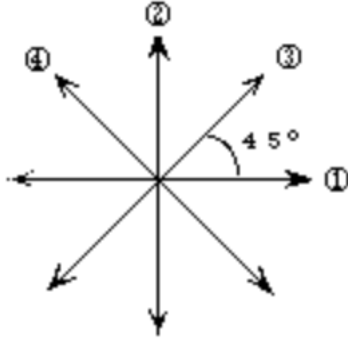
するエリート待望論を持っていて、その点でニーチェに共鳴した面もあります。当時のアレヴィ・サークルには、二〇世紀前半ロマン・ロランとともにフランスの左翼知識人運動の中心を担ったアンリ・バルビュスや、のちの人民戦線の首班レオン・ブルム、そしてマルセル・ブルーストがいました。また、ブルーストと同世代のジイドが、のちのフランスの世代にとって、ニーチェへの橋渡し役を務めることになりました。二〇世紀に入りますと、三〇年代のナチによるニーチェの横領が生じますが、それと思想的に真つ向から闘った人として、ニーチェの影響を強く受けつつ自己形成をしたジョルジュ・バタイユがあげられます。バタイユにとってニーチェとファシズムという問題は避けられない問いでした。逆に言う、三〇年代からニーチェとファシズムという問いは、フラ

OCRシステム

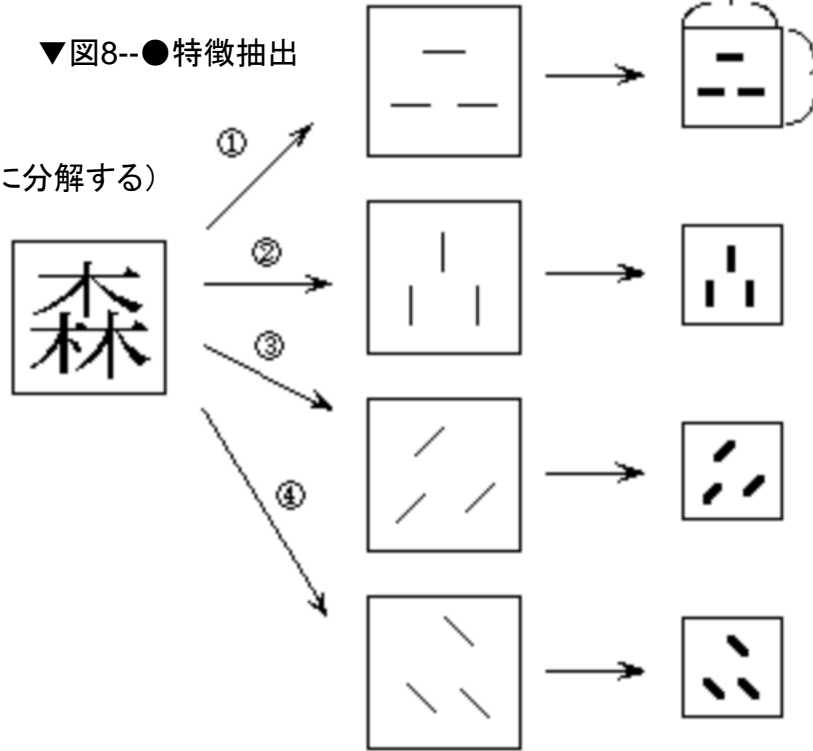
- 市販ソフト
 - メディアドライブ WinReaderPRO、 e.Typist
 - ABBYY FineReader
 - パナソニック 読取革命
 - . . .
- サービス
 - Google Cloud Vision API
 - Evernote API
 - . . .

OCRの仕組み

▼図9●4つの方向(1~4の方向に近い線分に分解する)



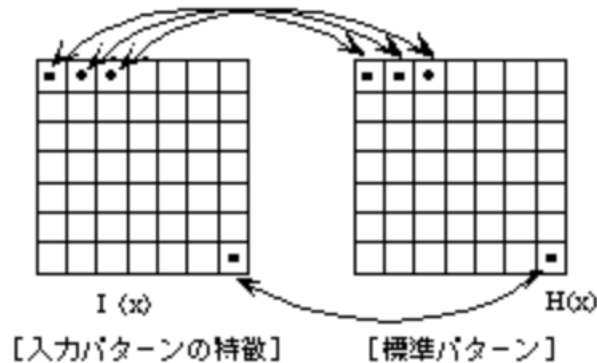
▼図8--●特徴抽出



(正規化パターン) → (方向成分に分解) → (圧縮)

的な特徴]

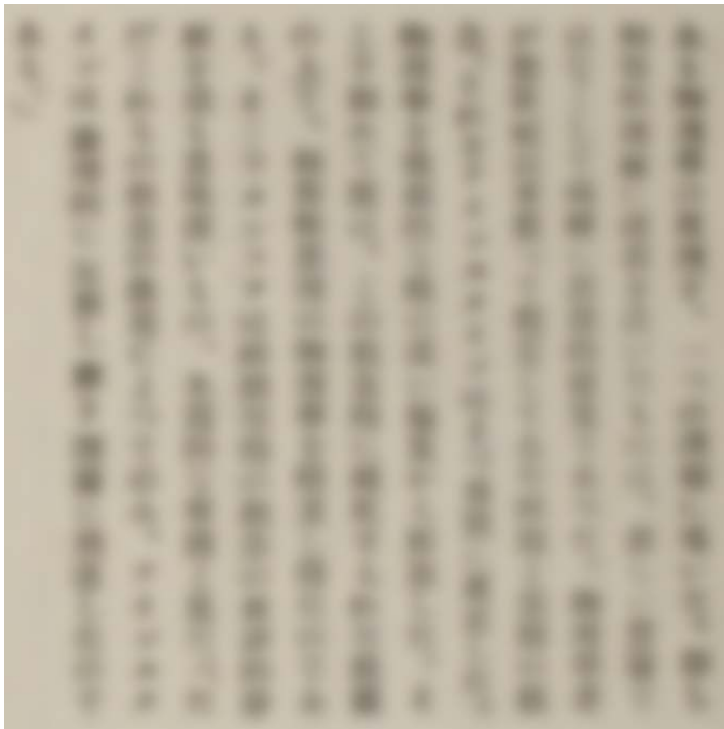
▼図10●ユークリッド距離(D)



出典:
<https://mediadrive.jp/technology/techocr05.html>

古い文献のデジタル化

- デジタルアーカイブの対象は古い文献が多い



ある物理学の原理を、一つの調和に導いた。即ち相対性理論に礎石を置いたものは、新しい実験ではなくして純粹し思想的探究であった。物理堅者が數世紀以來黙って暇定してゐた時間と空間の概念、それをアインスタインはまづ勇敢に更正した。物理学を傳統的な根の深い偏見から解放した。そこで初めて彼は、この概念的に純化せられた根據の上で、相對性原理の物理学を樹立し得たのである。そこでカシラアは時間空間の概念の批評的分析を最も意深きもの、本質的な契機と見た。「ただこれらの概念の改造によつてのみ、アインスタインは論理的に反對し難き理論に到達したのである。」

古い文献では精度が低下

『思想』創刊号、1921年
MM「世界見聞録」p82より引用

OCR精度の向上による 言語処理結果の向上

- OCR精度

95% から 99.85% へ

OCR精度95%での用語抽出結果

	1930年	1931年	1932年	1933年	1934年	1935年	1936年	1937年	1938年	1939年
1	自然科学	日本	日本	道穂	日本	自己自身	西田哲学	道穂	アメリカ	自己同一
2	資本主義	カント	ギルヘルム	フランス	宗敦	自己限定	日本	自己自身	英国	日本
3	日本	アメリカ	ブルジ	自然科学	フランス	道穂	自己自身	日本	日本	道穂
4	フランス	自己自身	カント	カント	資本主義	日本	自己限定	アントニ	自己同一	プラトン
5	道穂	ブルジ	宗敦	日本	ブルジ	根附	寺田先生	主体	フランス	宗敦
6	ブルジ	自己限定	自然科学	ブルジ	昭和八	自然科学	フランス	職業生活	イギリス	ピュタゴラス
7	回顧	宗敦	ファウスト	アジア	道穂	フランス	自己矛盾	フランス	ヨーロッパ	自己自身
8	エンゲルス	フランス	資本主義	スペンサー	昭和九	独逸	ノエシス	表現活動	自己自身	支那
9	唯物弁証法	エンゲルス	自己自身	ドイツ	ドイツ	宗敦	道穂	ベルグソン	自由主義	絶対矛盾
10	宗敦	メーリング	絶対精神	ゲシヒテ	自然科学	口子	カント	自己否定	自由主義	吉利

道徳、宗教、ブルジョア等の認識間違いが多数みられる

精度低下の原因

- 異体字・旧字体

例: 教↔教 学↔學

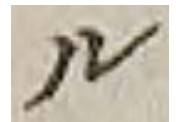
- フォントの違い



「と」:一画目が斜め、二画目の入りが大きい
→「ご」や「ざ」と誤る



「感」:「心」の部分が小さい
→「戚」や「咸」と誤る

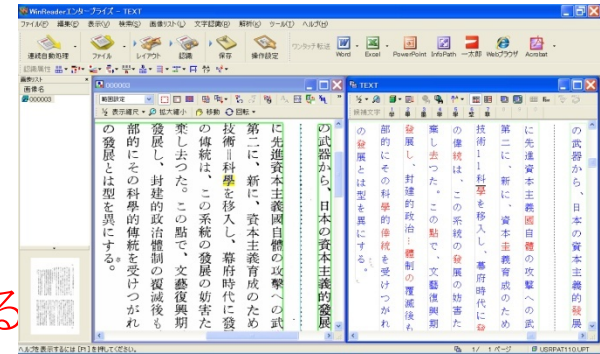


「ル」:右側が上がっている
→「ル」_ルと誤る

OCRによる電子化と課題

- 文字認識精度が低い
 - 旧字体の認識精度の問題
 - 強調部分、ルビ、外来語部分の認識
 - レイアウト解析に失敗し文字化けする

→近代文語論文の特徴抽出を行い精度を上げる



- 人手による作業コストの増大
 - レイアウト解析が不十分
 - タイトル、著者、脚注部分等のメタ情報を識別しなければならない
 - 論文の区切りを識別（ページの途中で別論文が始まる等）
 - OCRシステム開発会社と協同でOCRの高精度化
 - 誤認識の人手による訂正と学習
 - 異体字の対応、ルールによるレイアウト解析

→自動化・システム化を行い、作業コストを軽減する



精度99.85%での用語抽出結果

	1930年	1931年	1932年	1933年	1934年	1935年	1936年	1937年	1938年	1939年
1	自然科学	ヘーゲル	ヘーゲル	自然科学	ヘーゲル	自己自身	ヒューマニズム	アントニ	アメリカ	自己同一
2	マルクス主義	マルクス主義	ゲーテ	唯物論	マルクス	自己限定	西田哲学	ヘーゲル	ヘーゲル	ピュタゴラス
3	唯物論	弁証法	マルクス主義	マルクス	フランス	ヘーゲル	ヘーゲル	自己自身	自己同一	プラトン
4	ヘーゲル	カント	ギルヘルム	カント	キリスト	アリストテレス	自己自身	職業生活	自由主義	自己自身
5	エンゲルス	自己自身	形而上学	スペンサー	アリストテレス	自然科学	自己限定	表現活動	イギリス	ディオニューソス
6	プロレタリアート	プロレタリアート	弁証法	存在論	日本精神	ロゴス	弁証法	ベルグソン	自己自身	キリスト
7	資本主義	ジャーナリズム	エンゲルス	アリストテレス	自然科学	プラトン	ノエシス	自己否定	フランス	モラル
8	トーカー	形而上学	ファウスト	イデオロギー	エンゲルス	エネルギー	プラトン	プラトン	ヨーロッパ	ロゴス
9	コント	マルクス	ディルタイ	ディルタイ	自己自身	ノエシス	アリストテレス	弁証法	ロゴス	エートス
10	観念論	イデオロギー	カント	形而上学	カント	ベルグソン	ソクラテス	デカルト	アントニ	ソクラテス

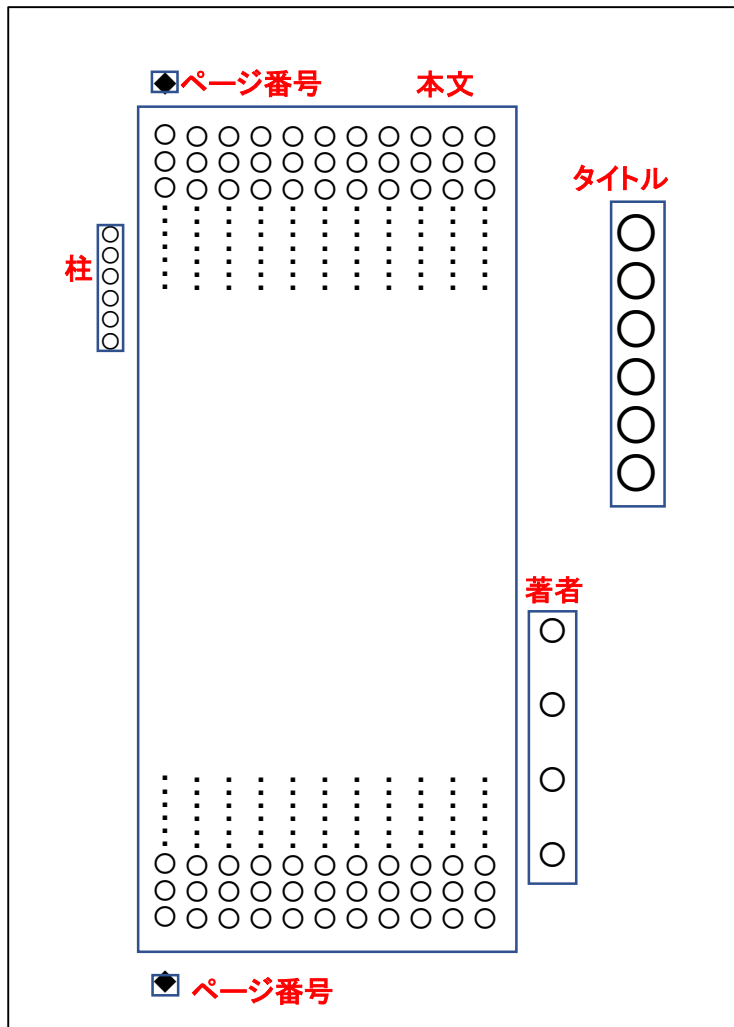
上位には間違いがほとんど見られない

「思想」の構造化プロジェクト

- テストとして20年分(1930-49)をテキスト化
 - OCRシステム開発会社と協同でOCRの高精度化
 - 異体字の対応、ルールによるレイアウト解析
 - 99%程度の精度でテキスト化
 - タイトル・著者等も自動的に抽出

→80年分に拡大

レイアウト構造の解析



- 各ブロックに論理ラベルを自動的に付与
- ルールベースで付与
 - 文字サイズ
 - 位置
 - ブロックの順番

レイアウト解析ソフトウェア の評価

- 3号分（250ページ、1631ブロック中）における解析間違い

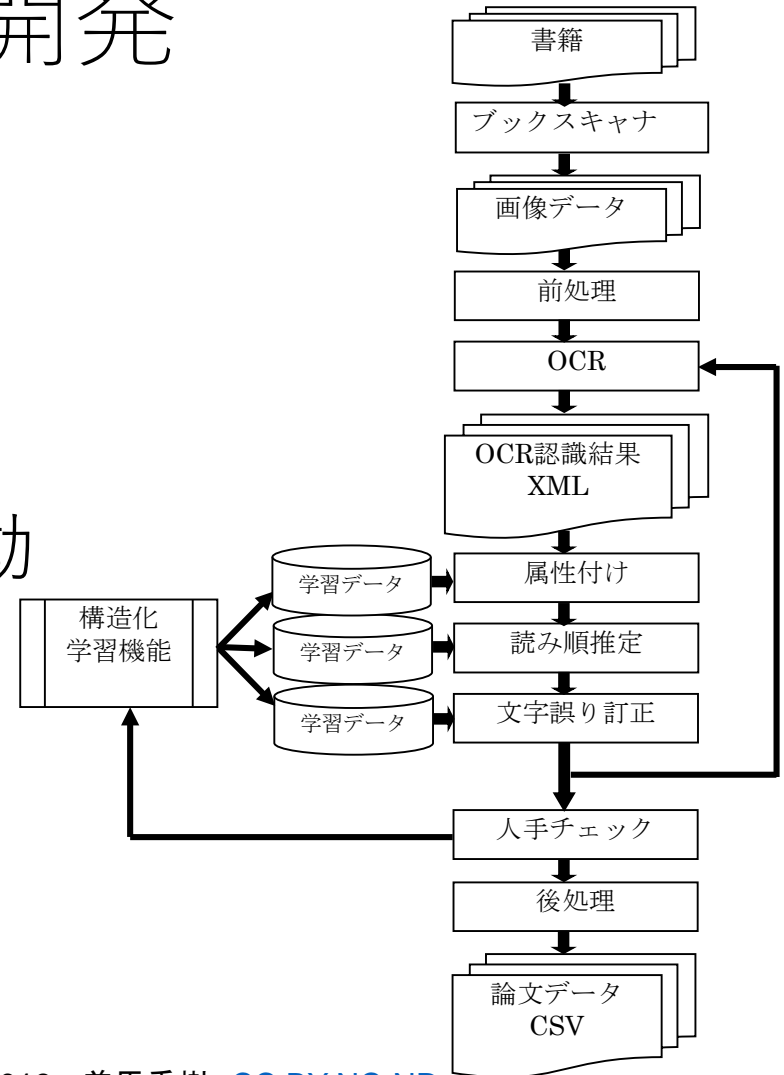
①ページ番号や欄外の見出しを本文の一部と誤認識する	51件
②属性付けの間違い（タイトル、著者）	2件
③レイアウト解析の間違い	4件
④ルビや句読点、汚れの存在によるレイアウト解析の間違い	5件
⑤広告や編集後記の存在による区切りの誤認識	5件
⑥特殊なレイアウトの存在による区切りの誤認識	3件

合計 70件

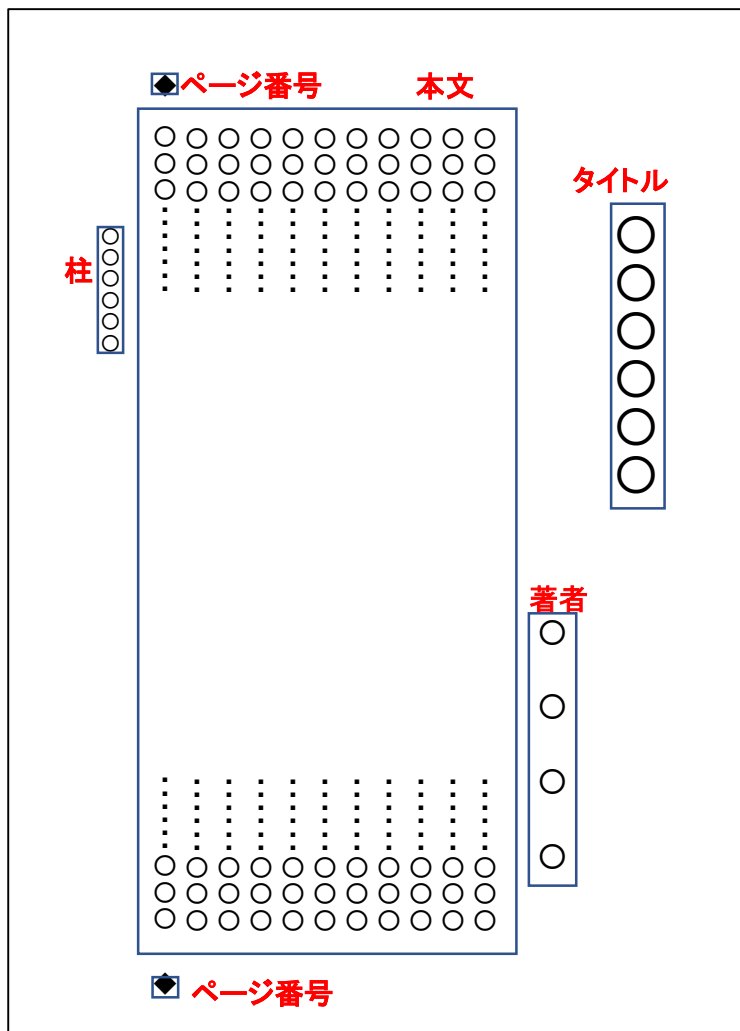
全1631ブロック中、1561ブロックを正しく認識（認識率95.7%）

現在の取り組み 人工知能（機械学習）を用いた 自動化システムの開発

- ルールベースの限界
 - 『思想』に特有のルールを追加すると...
 - ◎ 『思想』の認識精度は向上
 - × 他の文献には適用できない
- 機械学習を用いた新たな自動化システムの開発
 - 属性付け
 - 読み順推定
 - 統計的文字誤り訂正



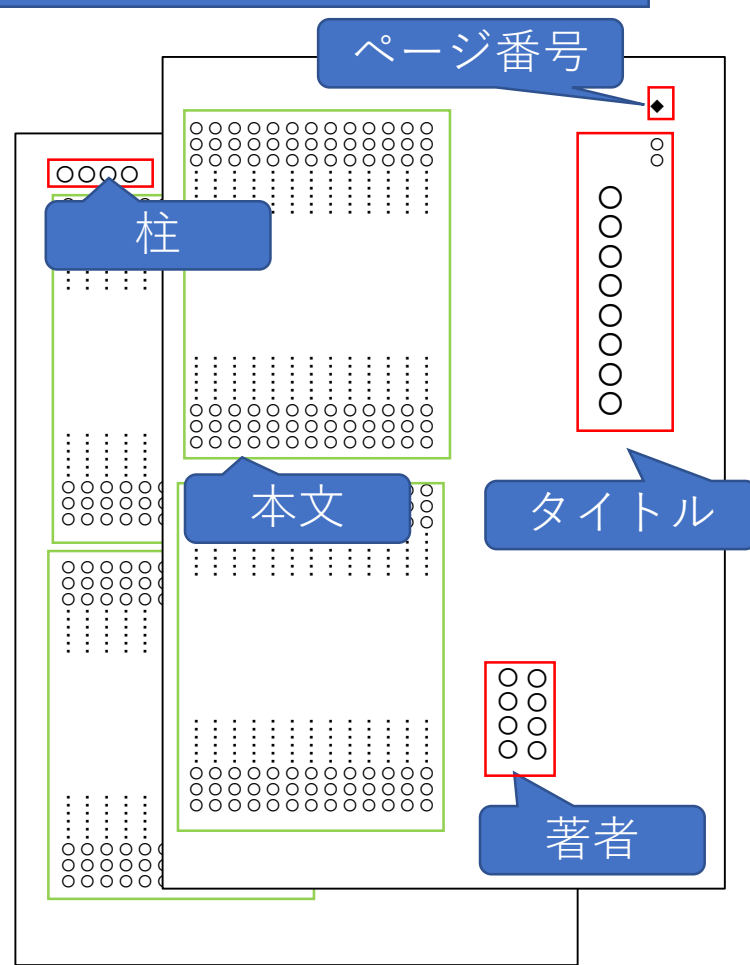
論理レイアウト構造の認識



- 各ブロックに論理ラベルを自動的に割り当てる
- 機械学習手法を利用
 - 様々なレイアウトに対応
 - 割り当てに使用する特徴
 - 文字サイズ
 - 位置
 - 周辺の余白
 - . . .

機械学習を用いた レイアウト属性付け

レイアウト属性の正解データ



特徴量による
分類器

特徴量

- 視覚的特徴量**
- 文字サイズ
 - ブロックサイズ
 - ブロック位置
 - 余白長さ
- 言語的特徴**
- 「名詞」割合
 - 「人名」割合

精度(F値)	ルール	機械学習
タイトル	0.913	0.960
著者	0.966	0.985
柱	0.974	0.979
ページ番号	0.994	0.995
本文	0.985	0.992

人手で作成したルールに比べ高精度

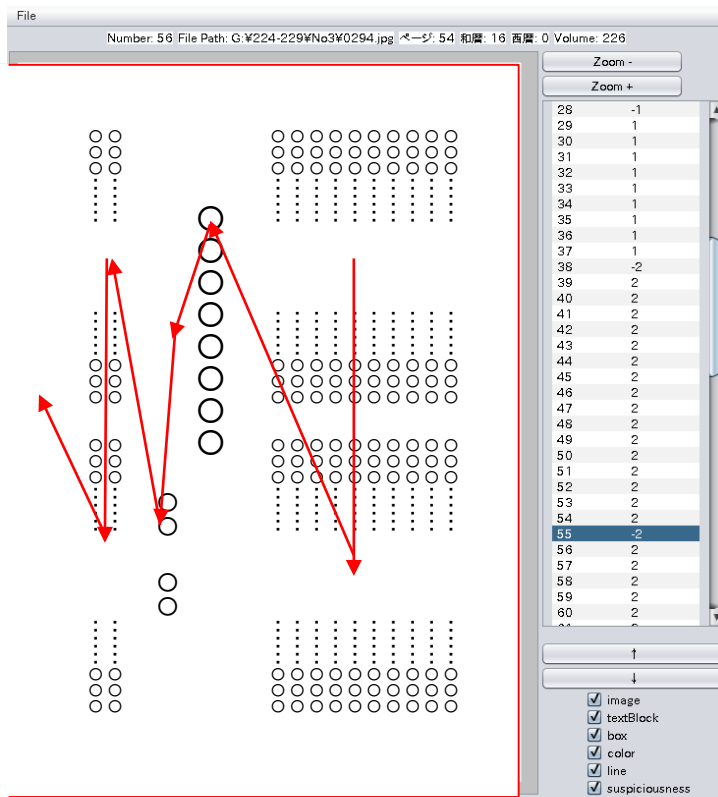
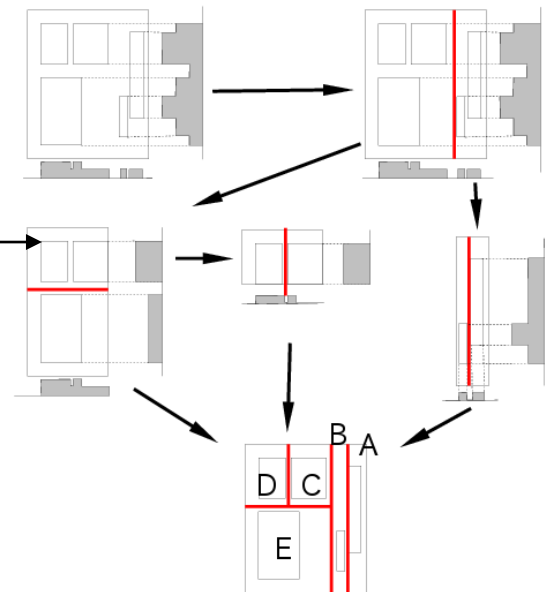
ブロック読み順の推定

OCR結果のブロック間の順番を自動で推定

自動読み上げ

テキスト検索の精度向上

機械学習による推定



人手による修正インターフェース

言語情報を用いたOCR誤りの訂正

OCR結果

心晴が真に心情的になるとき、自つから...

文字候補:
晴 惰 椿 情 蜻 清

文字候補:
づ つ っ

心晴が
心惰が
心椿が
心情が
心蜻が
心清が

もっともらしいのはどれか？

自つか
自づか
自っか

訂正結果

心情が真に心情的になるとき、自づから...

年代別文字認識精度

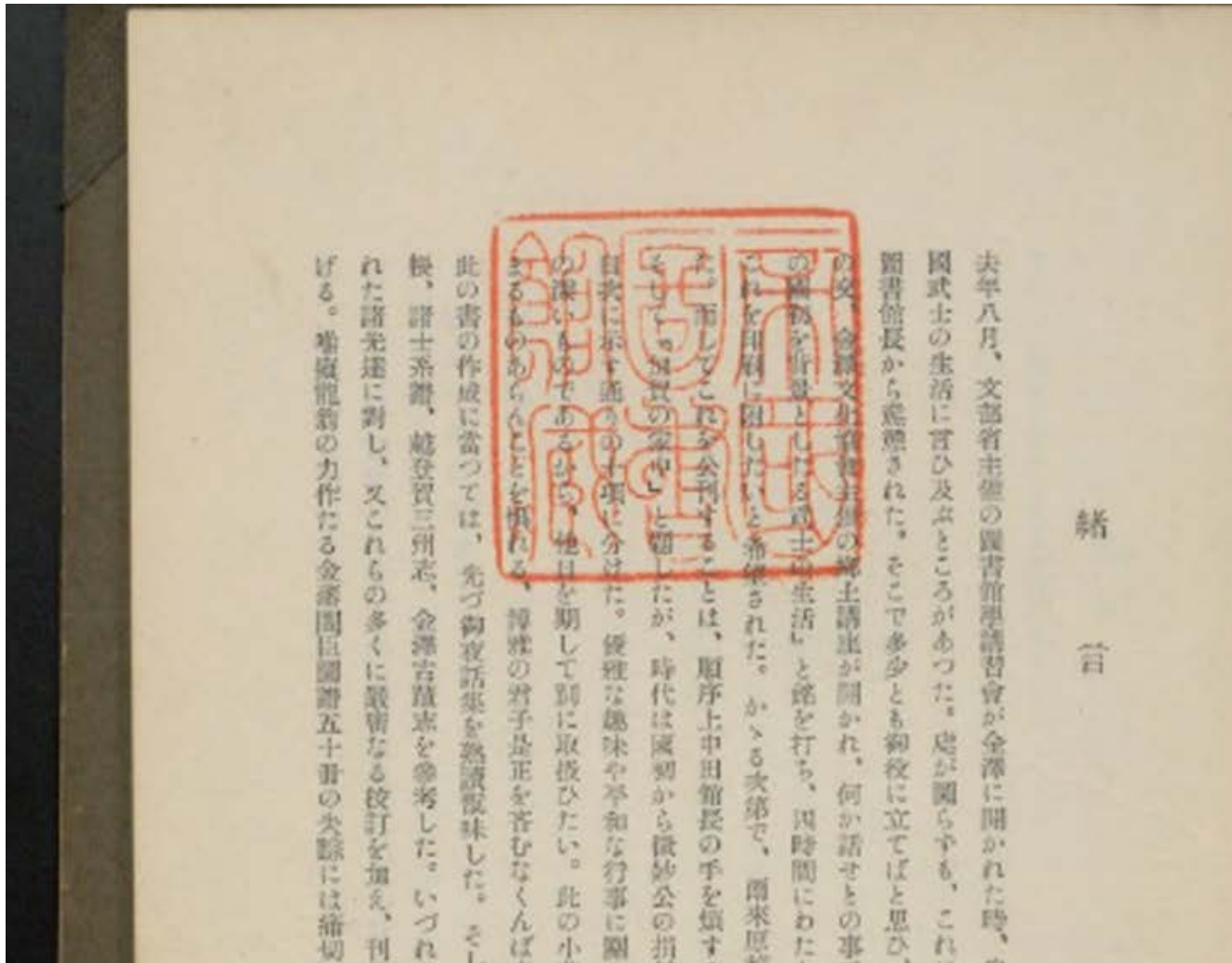
各年代1論文を人手でチェック

年	データ番号	文字数	誤認識数	認識精度
2000	0010-1	15998	31	99.81%
1990	0010-1	41638	215	99.48%
1980	0010-1	32229	83	99.74%
1970	0010-1	29476	846	97.13%
1960	0010-1	3090	61	98.02%
1950	0010-1	6066	170	97.19%
1940	0010-1	8316	52	99.37%
1930	0010-1	4294	348	91.90%
1921	0010-1	9218	1078	88.31%

古い年代で大きく精度が低下

実際に画像を見てみる と・・・

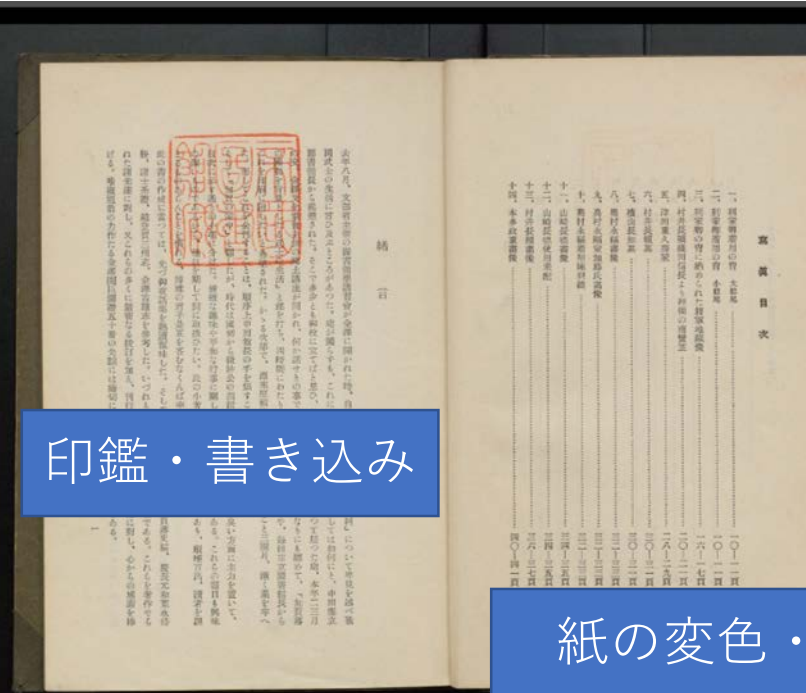
国立国会図書館デジタルコレクション
岡本勇著『加賀の家中』(昭和10年)



精度低下の原因

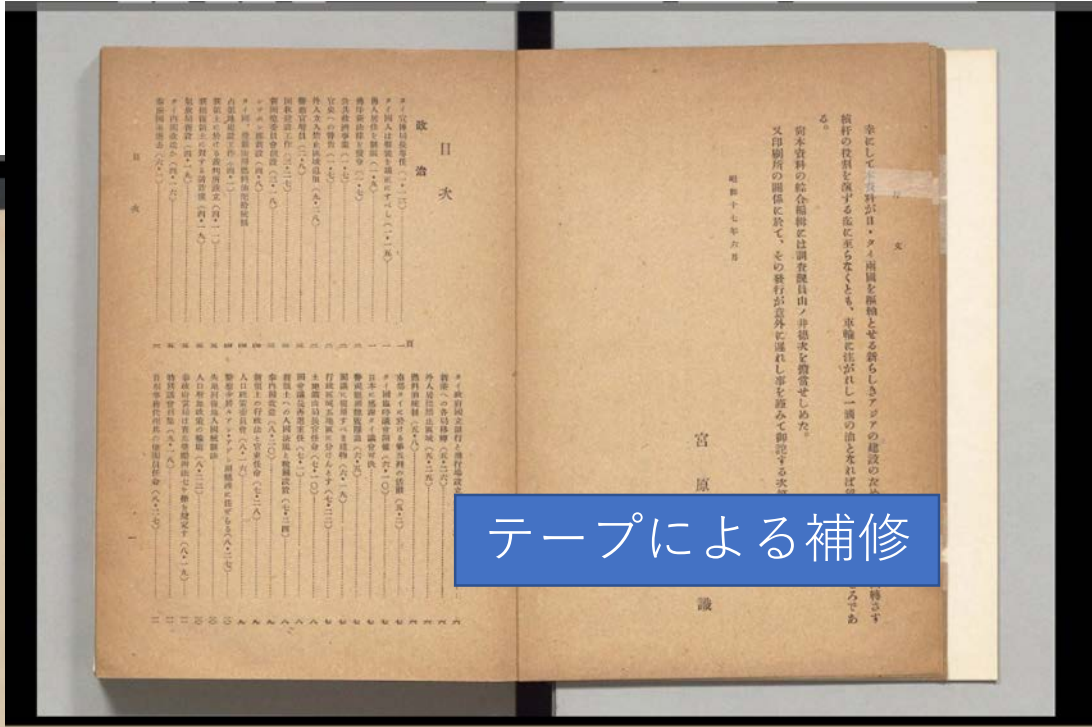
- 原本状態の悪さ

国立国会図書館デジタルコレクション
岡本勇著『加賀の家中』(昭和10年)



印鑑・書き込み

紙の変色・汚れ



テープによる補修

国立国会図書館デジタルコレクション
『一九四一年タイ国政治経済情勢』(1942年)
<http://dl.ndl.go.jp/info:ndljp/pid/1879315>

現状では対応が不可能

その他、OCRの課題

OCR処理: ページ端・背景の誤

認識

- ページの端や背景を文字と誤認識する
- 誤認識が他のブロックに影響する
 - ブロック境界誤り
 - 行順番誤り

ブ

SSEdit でのブロック再構成時に発生
→前処理でブロックを削除



番号	タイトル	著者名	開始頁	終了頁	論文頁	開始フ
1			6	18	3	0003_2
2	毒島の政治状況の存在を物		8	15	5	0003_2

20000010.xml - BSEdit

ファイル(E) 編集 解析 環境 表示(V) ヘルプ(H)

修 ページ

6
7
8
9
10
11
12
13
14
15
16
17
18
19

表示倍率
1/3

◆ ○○○○○○

○○○○

○○○..... ○○○○
○○○..... ○○○○
○○○..... ○○○○
○○○..... ○○○○

○○○○

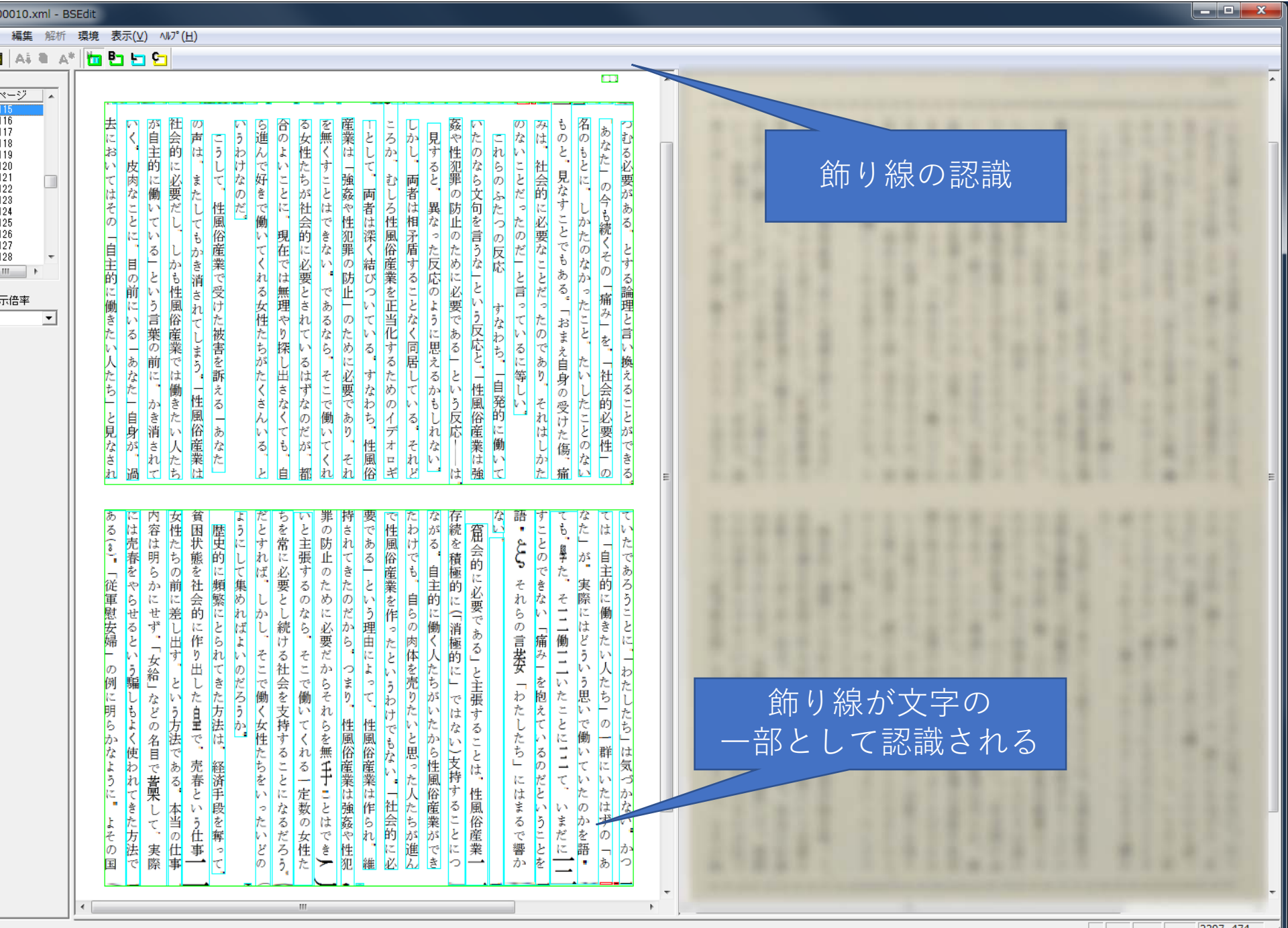
○○○..... ○○○○
○○○..... ○○○○
○○○..... ○○○○
○○○..... ○○○○

○○○○

○○○..... ○○○○
○○○..... ○○○○
○○○..... ○○○○
○○○..... ○○○○

○○○○

○○○..... ○○○○
○○○..... ○○○○
○○○..... ○○○○
○○○..... ○○○○



飾り線の認識

飾り線が文字の一部として認識される

つむむる必要がある。とする論理と言い換えることができる。あなた」の今も続くその「痛み」を、「社会的必要性」の名のもとに、しかたのなかったこと、たいしたことのないものと、見なすことでもある。「おまえ自身の受けた傷、痛みは、社会的に必要なことだったのであり、それはしかたのないことだったのだ」と言っているに等しい。

これらのふたつの反応 すなわち、「自発的に働いていたのなら文句を言うな」という反応と、「性風俗産業は強姦や性犯罪の防止のために必要である」という反応——は「見すると、異なった反応のように思えるかもしれない」。しかし、両者は相矛盾することなく同居している。それどころか、むしろ性風俗産業を正当化するためのイデオロギ「として、両者は深く結びついている。すなわち、性風俗産業は「強姦や性犯罪の防止」のために必要であり、それを無くすることはできない。であるなら、そこで働いてくれる女性たちが社会的に必要とされているはずなのだが、都合のよいことに、現在では無理やり探し出さなくても、自ら進んで好んで働いてくれる女性たちがたくさんいる。というわけなのだ。

こうして、性風俗産業で受けた被害を訴える「あなた」の声は、またしてもかき消されてしまう。「性風俗産業は社会的に必要だし、しかも性風俗産業では働きたい人たちが自主的に働いている」という言葉の前に、かき消されていく、皮肉なことに、目の前にいる「あなた」自身が、過去においてはその「自主的に働きたい人たち」と見なされ

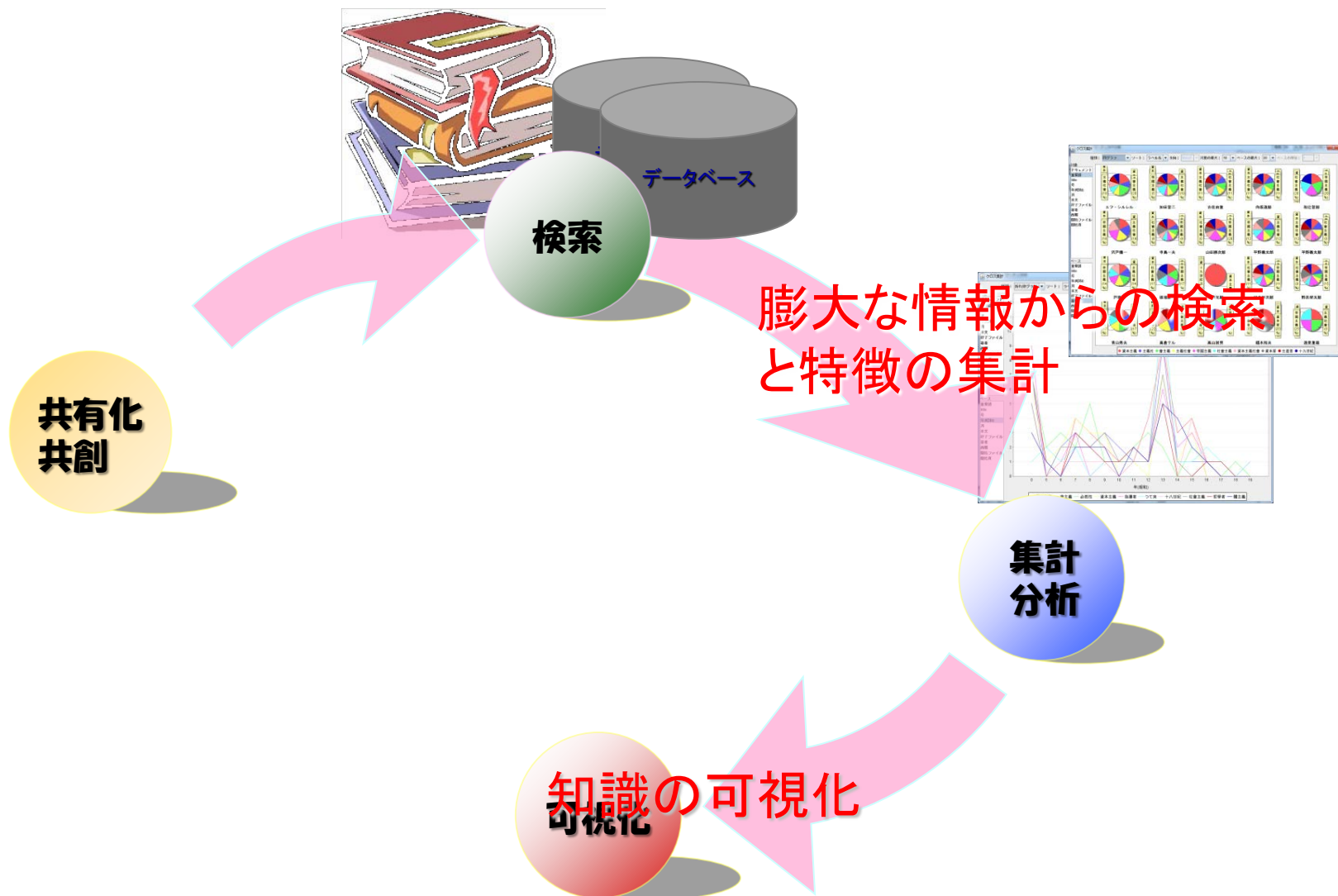
ていたであろうことに「わたしたち」は気づかない。かつては「自主的に働きたい人たち」の一群にいたはずの「あなた」が、実際にはどういふ思いで働いていたのかを語っても、果た、その「働いたこと」に「して、いまだに」すことのできない「痛み」を抱えているのだということに語・と、それらの言葉が「わたしたち」にはまるで響かない。

「社会的に必要である」と主張することは、性風俗産業「存続を積極的に（消極的に）ではない）支持することにつながる。自主的に働く人たちがいたから性風俗産業ができたわけでも、自らの肉体を売りたいと思った人たちが進んで性風俗産業を作ったというわけでもない。「社会的に必要である」という理由によって、性風俗産業は作られ、維持されてきたのだから、つまり、性風俗産業は強姦や性犯罪の防止のために必要だからそれらを無「くすることはできない」と主張するのなら、そこで働いてくれる一定数の女性たちを常に必要とし続ける社会を支持することになるだろう。だとすれば、しかし、そこで働く女性たちをいっただのようにして集めればよいのだろうか。

歴史的に頻繁にとられてきた方法は、経済手段を奪って「貧困状態を社会的に作り出した」自主で、売春という仕事「女性たちの前に差し出す」という方法である。本当の仕事内容は明らかにせず、「女給」などの名目で苦勞して、実際には売春をやらせるといふ騙しもよく使われてきた方法である。この「従軍慰安婦」の例に明らかかなように、よその国

MIMAサーチ「思想」 構造化デモ

「思想」の構造化のサイクル

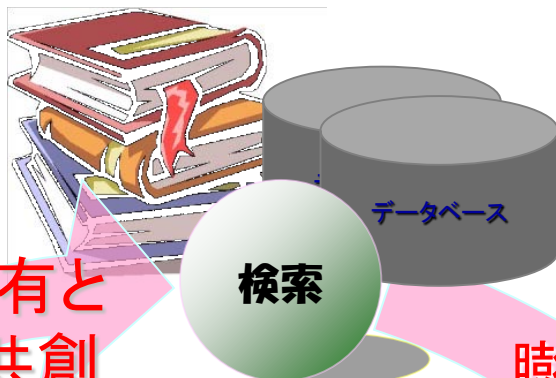


「思想」の構造化のサイクル



知識の共有と
再利用、共創

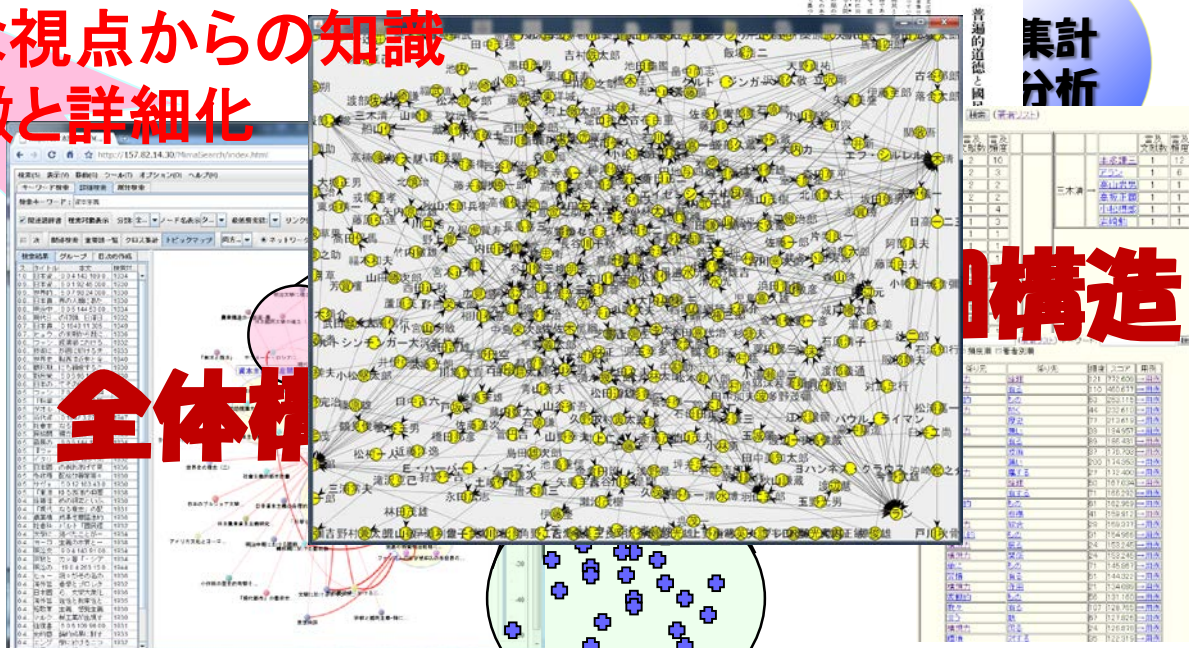
共有化
共創



膨大な情報からの検索
と特徴の集計



様々な視点からの知識
の俯瞰と詳細化



全体構

集計
分析

構造

「思想」構造化システム

- MIMAサーチ
- 係り受け関係検索
- 言及関係検索

MIMAサーチ

資本主義

総件数: 7451件中 1-50件目

S	タイトル	本文	検索
0.	現代論 第二次	19'	
0.	「ノイス 期を迎	19'	
0.	帝国主 ばじゅ	19'	
0.	現代に 一九六	19'	
0.	経済学 ないど	19'	
0.	経済学 ヒルファ	19'	
0.	帝国主 はじめ	19'	
0.	独占論 問題	19'	
0.	「資本 資本論	19'	
0.	スター 」するこ	19'	
0.	「資本 はじめ	19'	
0.	資本主 ないも	19'	
0.	帝国主 ここ数	19'	
0.	国家制 はじめ	19'	
0.	近代化 こんこ	19'	
0.	現代論 現代	19'	
0.	恐慌論 体系立	19'	
0.	国家、 国家独	19'	

著者

- 清水幾太郎 (68)
- 和辻哲郎 (47)
- 中村雄二郎 (43)
- 高山岩男 (42)
- 佐々木力 (39)
- 宮田光雄 (38)
- 三島憲一 (35)
- 広松渉 (32)
- 三木清 (31)
- 西田幾多郎 (31)
- 平田清明 (28)
- 宇野弘藏 (27)
- 日高六郎 (27)
- 杉山光信 (27)
- 井筒俊彦 (25)
- 南博 (25)

資本主義

マルクス

アメリカ

労働者

MIMAサーチ 詳細画面

Title 如何にして科学の発展は可能であるか

92号(1930年1月) p. 25

Author **三木清**

1/20

科学の発展に関して普通に行はれてある見解は次の如くである、科学は「進歩」(Fortschritt)の過程にある。しかもそれは不斷に、連続的に進歩してある。そしてこの点に於て科学は他の種々な文化形態に対して特殊な位置を占めてある。例へば芸術については、或る時代の芸術作品が初めて新しい技術的手段または遠近法の法則を創るに到つたからといって、そのためにそれが、それ以前の、かかる手段及び法則の知識を全く欠いてゐた芸術作品よりも、純粹に芸術的に一層高く立つてゐるとは云はれない。各の芸術作品はそれぞれの「完成」(Rundum)である、従つてそれは決して打ち克たれるものでなく、また決して旧びるものでもない

[もっと読む](#)

係り受け関係検索

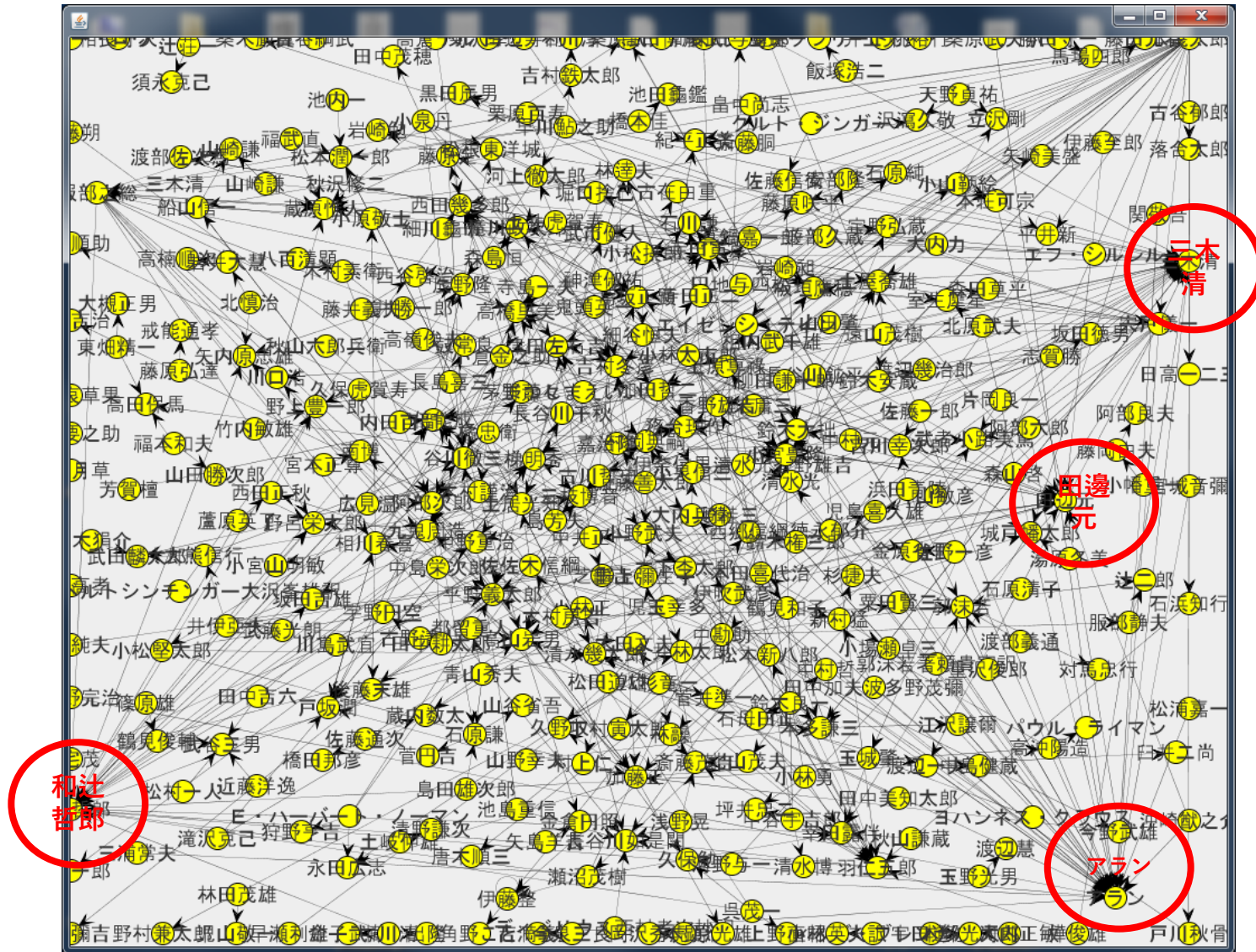
著者: (「思想」著者リスト) キーワード: 開始年: 終了年:

係り受け単位 キーワード単位

係り受け単位での検索結果
 (→キーワード単位で著者:「三木清」で検索)

著者	係り元	係り先	頻度	スコア	用例	三木清 著作リスト
三木清	構想力	論理	118	171.84	用例	ボルツァーの「命題自体」(1923/12 26号)
	言う	こと	708	143.33	用例	消息一通(1924/3 29号)
	こと	出来る	630	139.72	用例	パスカルと生の存在論的解釈(1925/5 43号)
	此の	もの	137	132.25	用例	愛の情念に関する説—パスカル覚書(1925/8 46号)
	もの	もの	332	124.81	用例	パスカルの方法(1925/11 49号)
	意味	於く	185	121.60	用例	パスカルの「賭」(1925/12 50号)
	もの	為る	476	120.62	用例	解釈学的現象学の基礎概念(1927/1 63号)
	こと	因る	281	116.69	用例	人間学のマルクスの形態(1927/6 68号)
	主観的	もの	69	111.39	用例	マルクス主義と唯物論(1927/8 70号)
	意味	有する	71	103.86	用例	プラグマチズムとマルキシズムの哲学(1927/12 74号)
	考える	こと	108	100.02	用例	ヘーゲルの歴史哲学(1929/4 83号)
	我々	出来る	70	99.41	用例	社会と自然(1929/8 87号)
	此の	場合	100	98.56	用例	紹介と検討 現象学叙説—山内氏の著作の読後記(1929/9 88号)
	我々	取る	56	98.20	用例	弁証法に於ける自由と必然(1929/10 89号)
	存在	仕方	92	96.87	用例	如何にして科学の発展は可能であるか(1930/1 92号)
	ところ	もの	100	96.59	用例	存在の歴史性(1931/12 115号)
	斯かる	もの	64	96.24	用例	拙著批評に答ふ(1932/9 124号)
	これ	反する	67	95.29	用例	倫理と人間(1933/6 133号)
	構想力	於く	41	94.95	用例	ヘルデル大辞典について(1934/9 148号)
	於く	言う	120	94.51	用例	表現に於ける真理(1935/9 160号)
	客観的	もの	66	93.05	用例	西田哲学の性格について(1936/1 164号)
	こと	無い	191	91.79	用例	ヒューマンイズムの哲学的基礎(1936/10 173号)
	もの	考える	139	88.93	用例	ヒューマンイズムの哲学的基礎(完)(1936/11 174号)
	もの	有る	237	88.59	用例	神話—構想力の論理に就いて(其一)(1937/5 180号)
	於く	為る	133	87.52	用例	神話(中)構想力の論理に就いて其二(1937/6 181号)
						制度(一)—構想力の論理に就いて(其四)(1937/8 183号)
						制度(二)—構想力の論理に就いて(其五)(1937/9 184号)
						制度(三)—構想力の論理に就いて(其六)(1937/10 185号)
						技術—構想力の論理に就いて其七(一)(1938/2 189号)
						技術—構想力の論理に就いて其八(二)(1938/3 190号)
					技術—構想力の論理に就いて其九(三)(1938/5 192号)	
					経験(一)—構想力の論理に就いて 統一(1939/9 208号)	
					経験(二)—構想力の論理に就いて 統一(1940/8 219号)	
					経験(三)—構想力の論理に就いて(統一)(1940/11 222号)	
					経験(四)—構想力の論理に就いて(統一)(1940/12 223号)	
					経験(五)—構想力の論理に就いて(統一)(1941/1 224号)	
					経験—構想力の論理に就いて(統一)(1941/2 225号)	
					充足的経験論(1941/4 227号)	
					経験(七)—構想力の論理に就いて(統一)(1942/4 239号)	
					経験—構想力の論理に就いて(統一)(1943/3 250号)	
					経験—構想力の論理に就いて(統一)(1943/11 251号)	

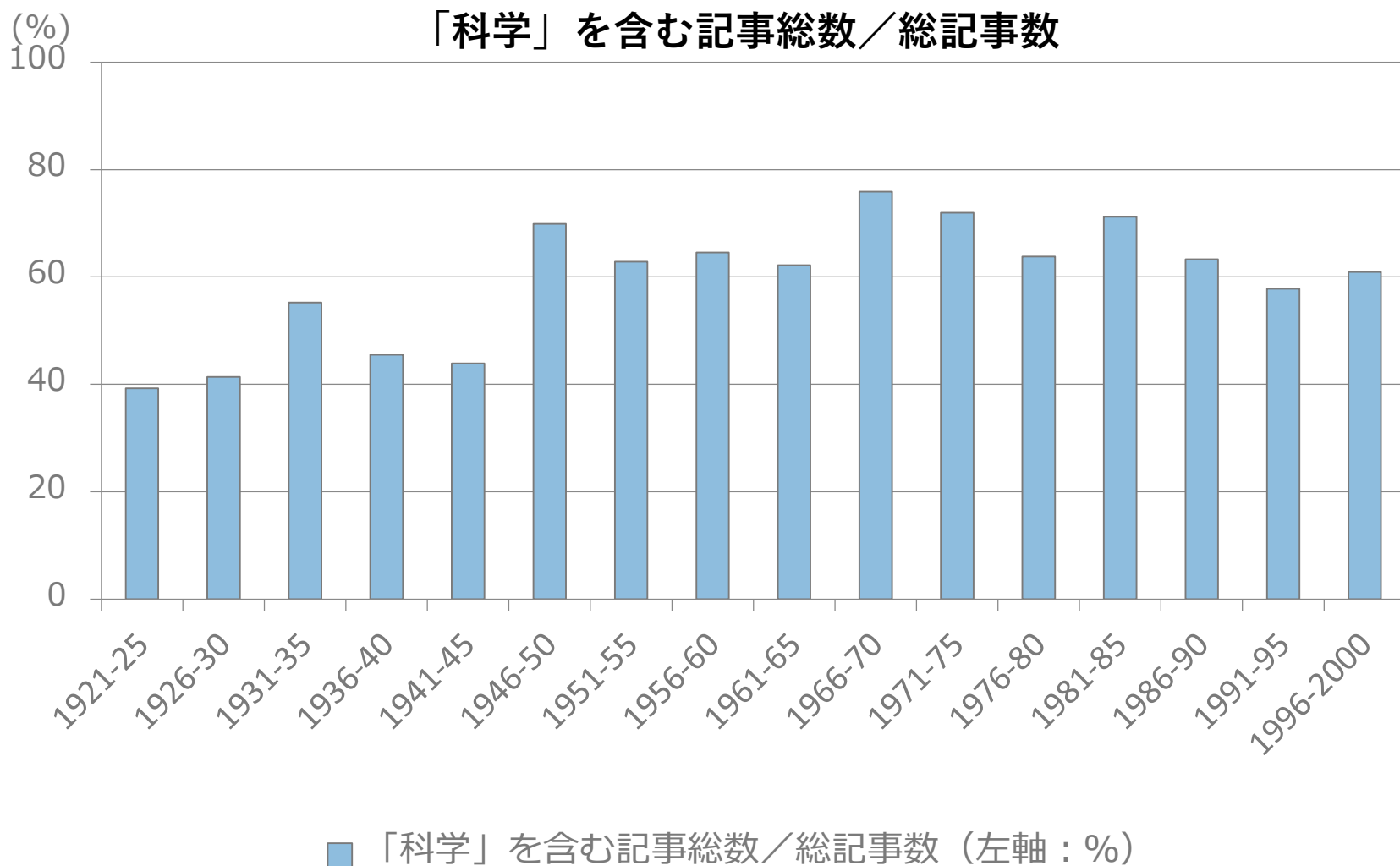
著者間の言及関係の可視化



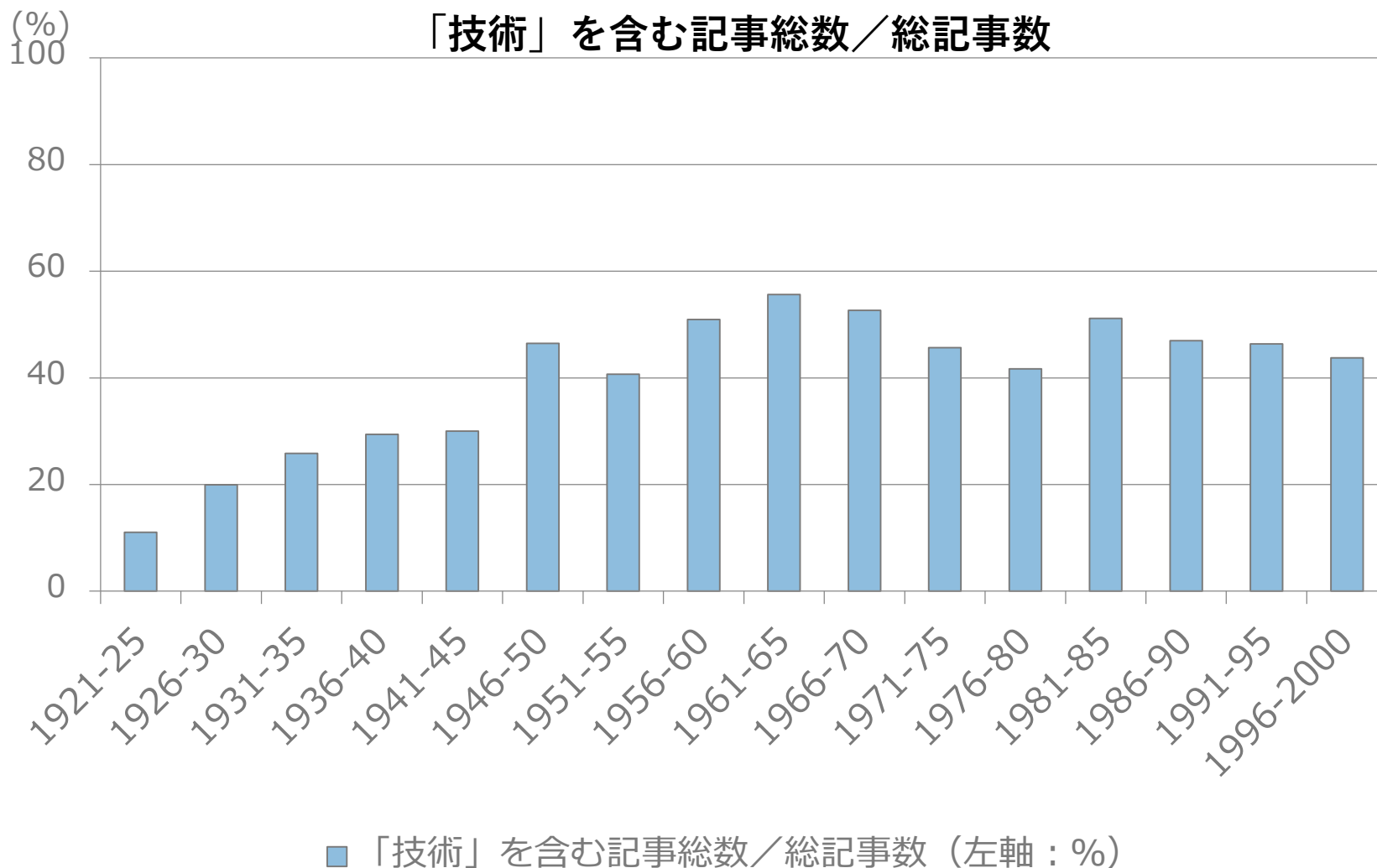
『思想』の分析例

- 使用される言葉・概念の年代による変遷
 - 対象1: 「科学」「技術」「科学技術」
 - 似た概念の言葉の使われ方の違い
 - 対象2: 原子力の戦争利用と平和利用
 - 戦争利用: 「核兵器」「原子爆弾」「核実験」等
 - 平和利用: 「原子力発電」「原発」等

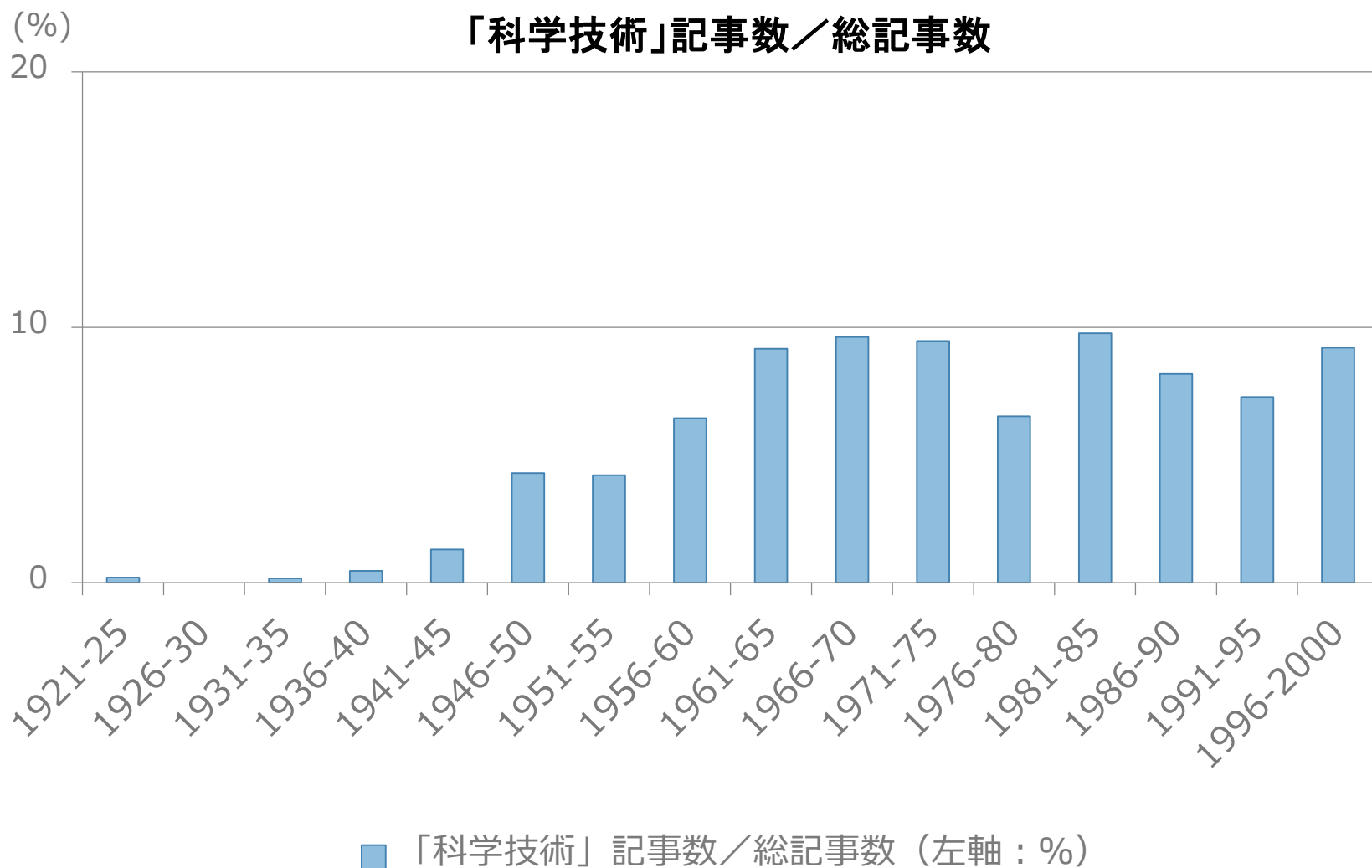
科学の登場頻度の推移



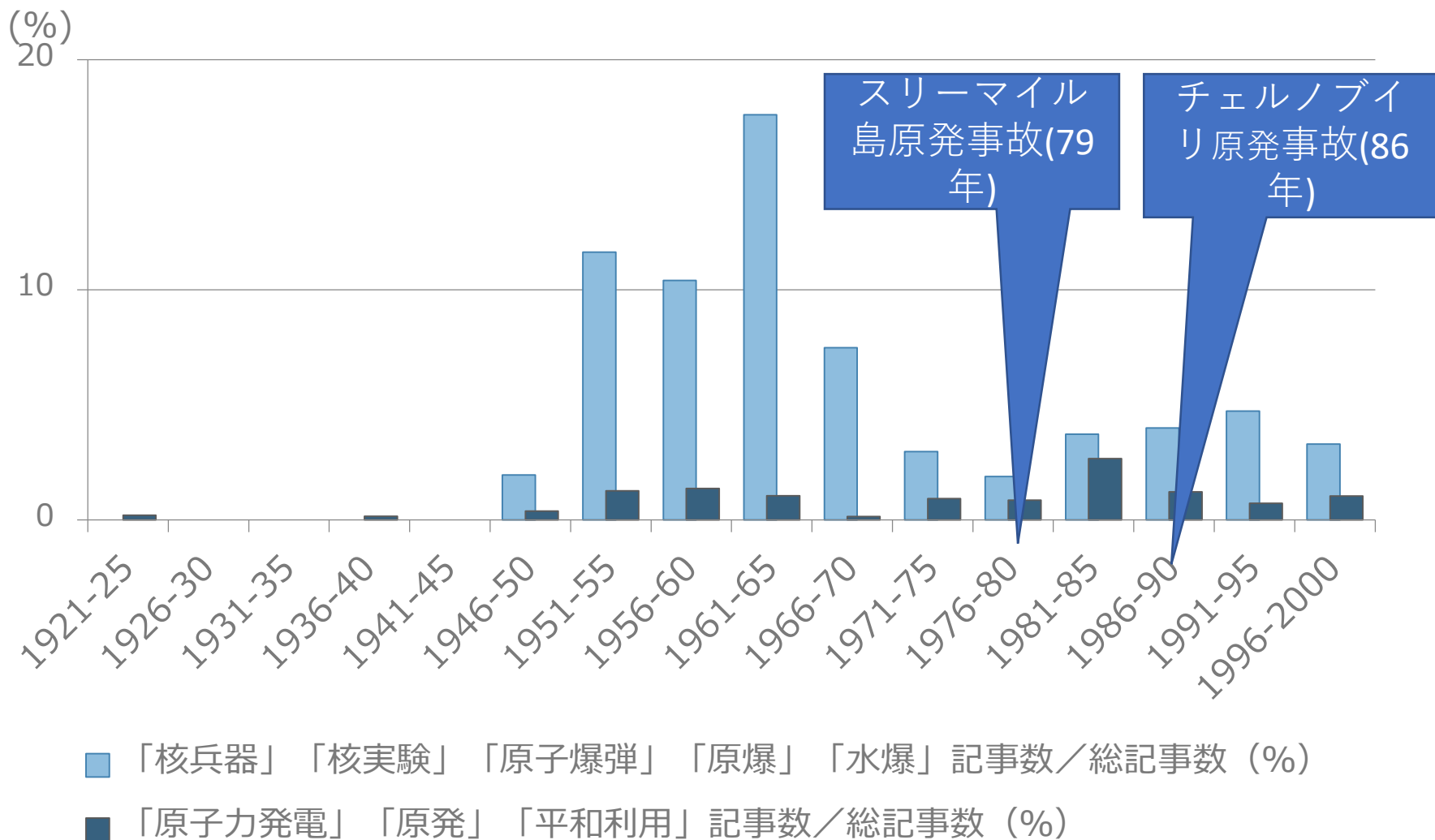
技術の登場頻度の推移



“科学技術”の登場頻度の推移



核兵器・原子力



講義／ワークショップの様子

世論調査 文学・思想・雑誌 政治思想史 読書・出版文化 論壇・メディア

出版流通

ナショナリズム

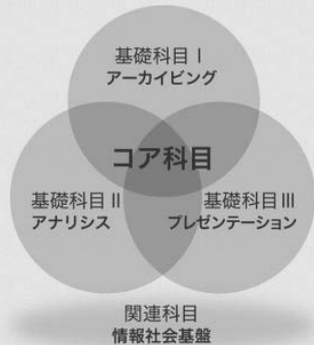
社会意識



技術開発チームとの対話

プログラム構成

本プログラムは、〈コア科目〉、〈基礎科目〉、〈関連科目〉によって構成されます。〈コア科目〉はデジタル・ヒューマニティーズの中核をなすもので、2つの「必修科目」が含まれます。〈基礎科目〉は、I)アーカイビング、II)アナリシス、III)プレゼンテーションという3つの要素からなり、領域を横断して理論と方法を学ぶことができます。〈関連科目〉では、これらの科目に関係する情報社会基盤の知識を得ることができます。



大学院教育への展開

- デジタルヒューマニティーズ
知の構造化論

コア科目とは



デジタル・ヒューマニティーズの中核をなす科目です

知の構造化論

講義題目	開講研究科	科目番号	担当者	単位	時間割	教室割
知識と情報処理、 人文学とコンピュータ、 臨床医療とコンピュータ	学際情報学府	4990110	美馬 秀樹 荒牧 英治	2.0	夏・水・4限 (14:50~16:20)	福武ホールB2階・ 福武ラーニングスタジオ1

◆ 授業の目標・概要 (クリックで詳細を表示)

本講義では、大量に蓄積された知識に対し、様々な構造化技術を用いて価値の創出を行う「知の構造化」を実践することを目的とする。知の構造化では、単なる検索にとどまらず、コンピュータを使って知識の抽出と分析を行い、分野や時勢を越えて各要素の間の関係を明らかにすることでその活用を促す。そのために、自然言語処理、人工知能、Web工学等の最先端技術を駆使する手法を学び、実社会で生じる様々な課題に対する分野を超えた知の適用の議論を進める。これにより、構造の抽出と、構造から機能を導出するプロセスを体現する。

人文学オープンデータ 共同利用センター

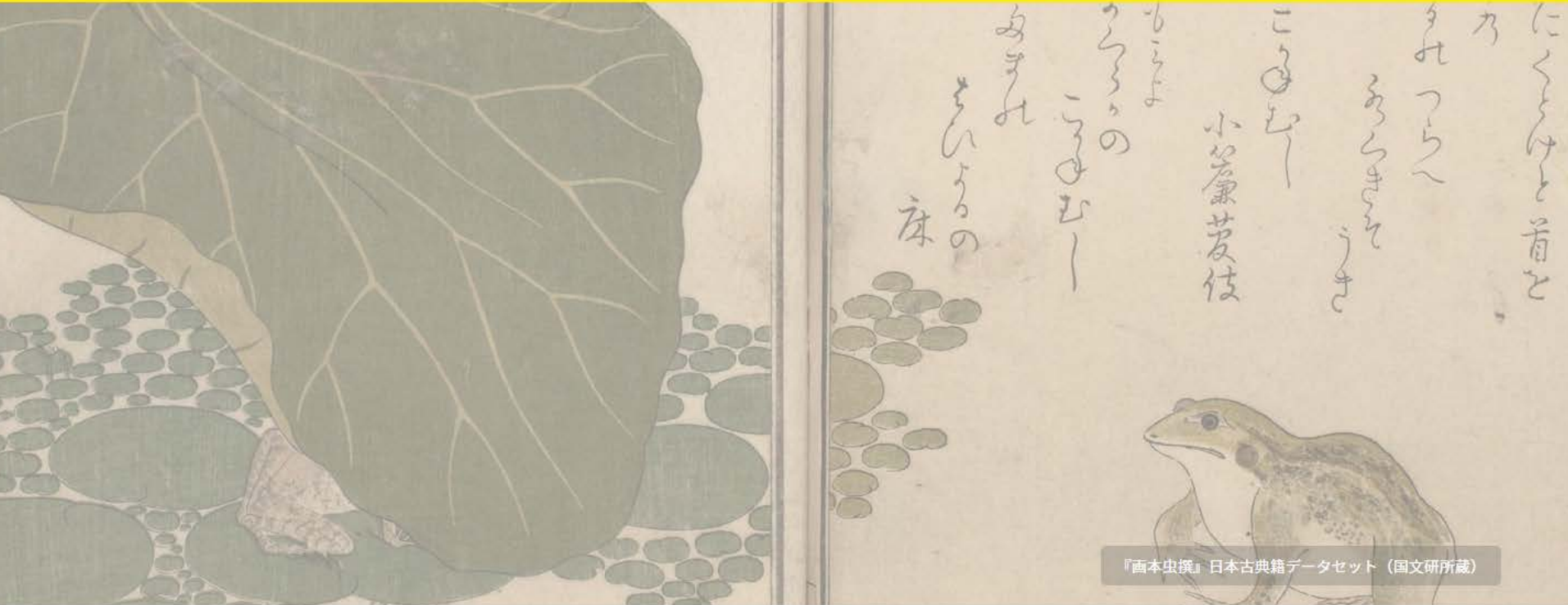
<https://codh.rois.ac.jp/>

CC BY-SA 4.0

 人文学オープンデータ共同利用センター
Center for Open Data in the Humanities

日本語 / English

メニュー



人文学オープンデータ共同利用センター (Center for Open Data in the Humanities / CODH) は、情報学・統計学の最新技術を用いて人文学データへのアクセスを改善する研究開発を進めるとともに、オープンサイエンスの考え方に基づき多くの人々が参加できるデータプラットフォームを構築することで、データ駆動型の人文学研究や超学際的な人文学研究など、情報技術を用いた新しい人文学の方法論を開拓します。[もっと詳しく..]

お知らせ

まとめ

- 文理融合による教育／研究の可能性
 - 可視化、対話による共創、知見の拡大、共有
- 資料体へのアクセス・公開
 - 教育支援、教材づくり
 - 海外の研究者の支援
 - 近代の言語資源の価値化