

1. 自然言語処理とは

1.1 自然言語処理概略

▶ 言語

- ▶ プログラミング言語(人工言語)
- ▶ 自然言語(ex. 日本語、英語、中国語、...)

▶ 自然言語処理の目的

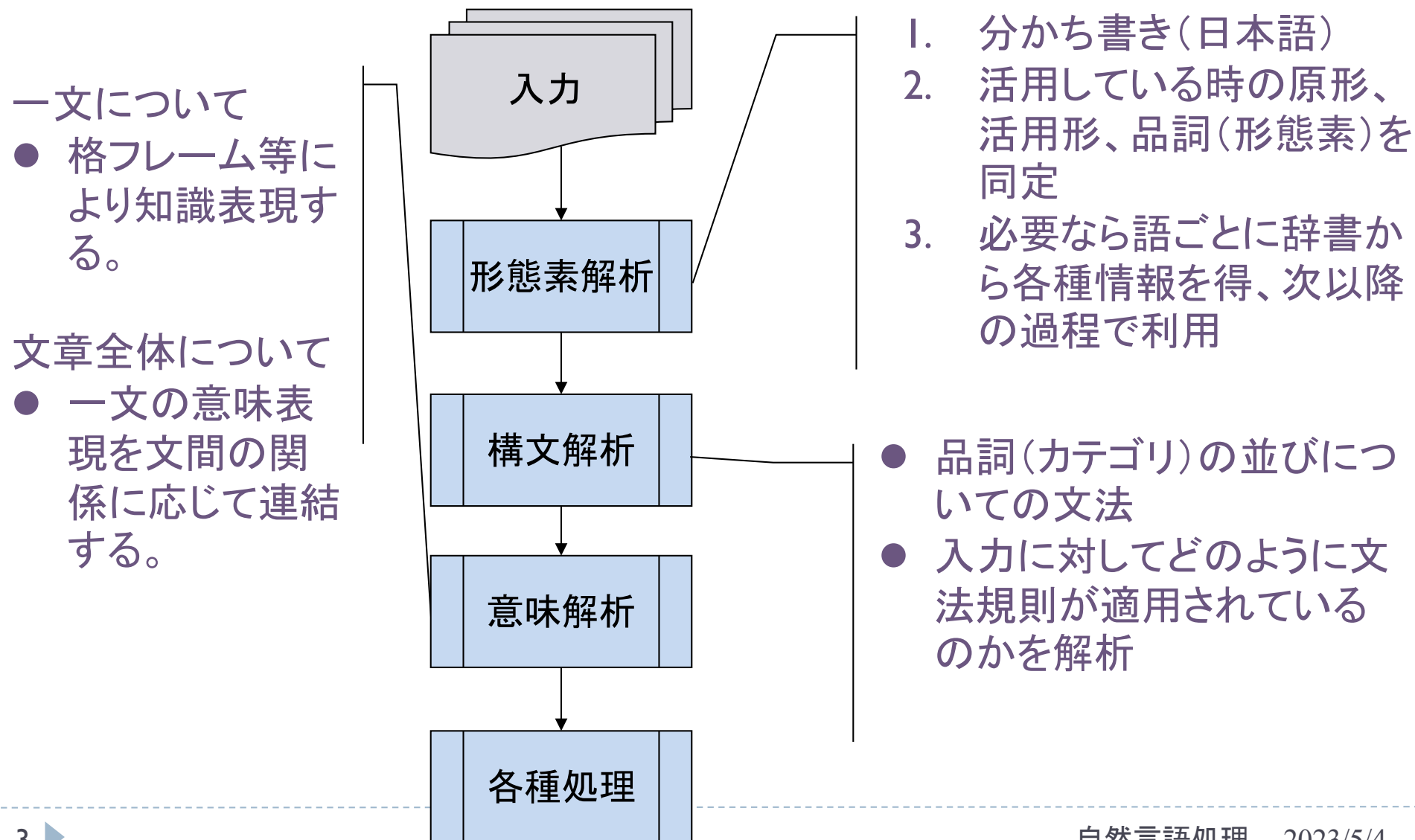
人間が使う言語を処理する

- ▶ 話す言葉、音声言語
- ▶ 書いた文,文章の解析
- ▶ 文章からの情報抽出
- ▶ 文,文章の生成

▶ 方法

- ▶ 認知的原理、モデル、知見に基づくものもあるが。
- ▶ 必ずしも認知的である必要はない。

自然言語処理過程の典型的モデル



文字の扱い ASCII Code Table

例 Book

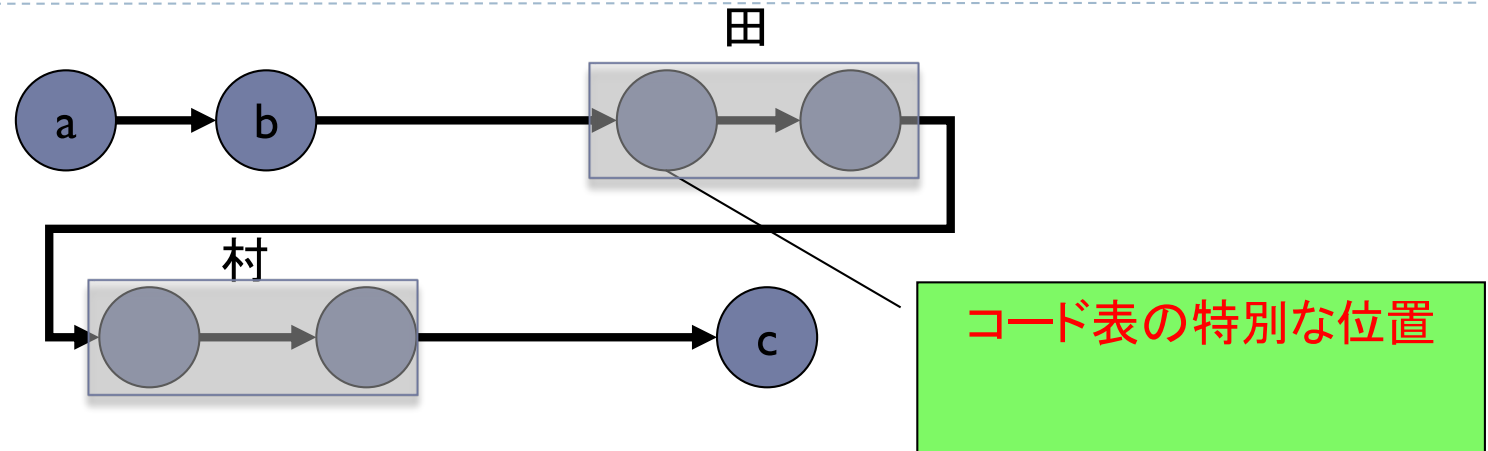
文字	二進数	16 進数
'B'	1000010	42
'o'	1101111	6f
'o'	1101111	6f
'k'	1101011	6b

b6	0	0	0	0	1	1	1	1
b5	0	0	1	1	0	0	1	1
b4	0	1	0	1	0	1	0	1
b3-b0								
0000	<i>nul</i>	<i>dle</i>			0	@	P	`
0001	<i>soh</i>	<i>dc1</i>	!		1	A	Q	a
0010	<i>stx</i>	<i>dc2</i>	"		2	B	R	b
0011	<i>etx</i>	<i>dc3</i>	#		3	C	S	c
0100	<i>eot</i>	<i>dc4</i>	\$		4	D	T	d
0101	<i>enq</i>	<i>nak</i>	%		5	E	U	e
0110	<i>ack</i>	<i>syn</i>	&		6	F	V	f
0111	<i>bel</i>	<i>etb</i>	'		7	G	W	g
1000	<i>bs</i>	<i>can</i>	(8	H	X	h
1001	<i>ht</i>	<i>em</i>)		9	I	Y	i
1010	<i>lf</i>	<i>sub</i>	*		:	J	Z	j
1011	<i>vt</i>	<i>esc</i>	+		;	K	[k
1100	<i>ff</i>	<i>fs</i>	,		<	L	¥	l
1101	<i>cr</i>	<i>gs</i>	-		=	M]	m
1110	<i>so</i>	<i>rs</i>	.		>	N	^	n
1111	<i>si</i>	<i>us</i>	/		?	O	_	o
								<i>del</i>

コード化 (EUC, SJIS, JIS)

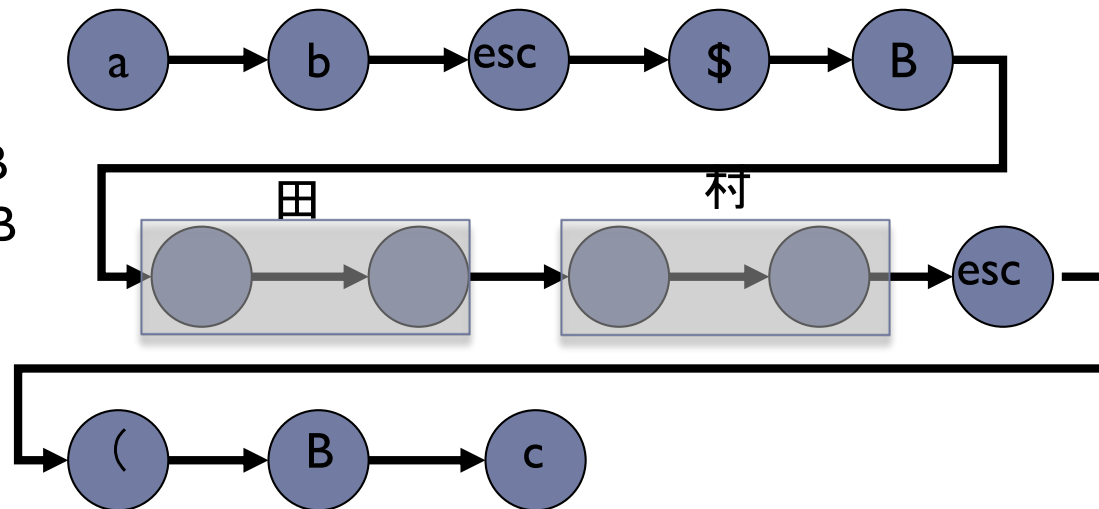
「ab田村c」の各コード

ShiftJIS

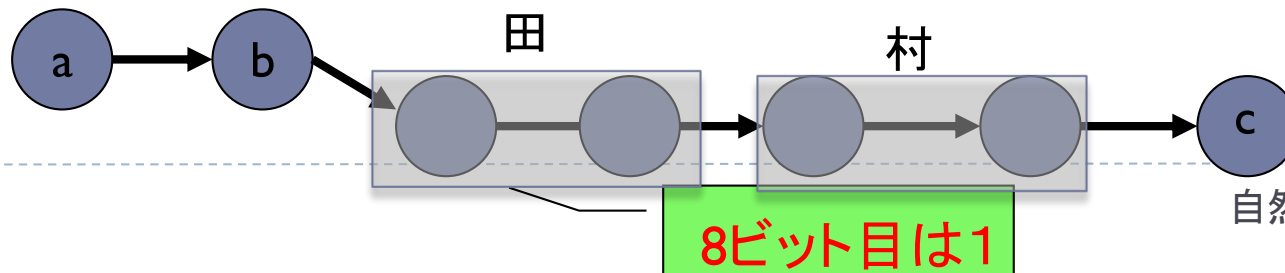


JIS

モード切り替え
漢字in: esc+\$+B
漢字out: esc+(+B



EUC



文字集合の例 (JIS X 0208)

文字集合(JIS X 0208)

- コンピュータで扱える文字の集合(Unicode)とその符号化(UTF-8)を
- 分離して考える。
- 一つの文字を区(row.1~94)点(cell.1~94)で表す。
例:「医」=16区69点(または16-69)
- 8,836(=94×94)字を扱うことができる。
- この構造は中国のGB2312、韓国のKS X 1001でも採用。

Row 16	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19
00		亜	啞	娃	阿	哀	愛	挨	始	逢	葵	茜	穉	惡	握	渥	旭	葦	鯨
20	梓	压	幹	扱	宛	姐	虻	飴	絢	綾	鮎	或	栗	裕	安	庵	按	暗	案
40	鞍	杏	以	伊	位	依	偉	圉	夷	委	威	尉	惟	意	慰	易	椅	為	畏
60	移	維	緯	胃	萎	衣	謂	違	遺	医	井	亥	域	育	郁	磯	一	壺	溢
80	稻	茨	芋	鰯	允	印	咽	員	因	姻	引	飲	淫	胤	蔭				

Row 17	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	
00		院	陰	隱	韻	吋	右	宇	烏	羽	迂	雨	卯	鵜	窺	丑	確	臼	渦	噓
20	唄	鬱	蔚	鰻	姥	厖	浦	瓜	閏	噂	云	運	雲	荏	餌	叡	當	嬰	影	映
40	曳	榮	永	泳	洩	瑛	盈	穎	穎	英	衛	詠	銳	液	疫	益	馭	悅	謁	越
60	閱	榎	厭	円	園	堰	奄	宴	延	怨	掩	援	沿	演	炎	焰	煙	燕	猿	緣
80	艷	苑	蘭	遠	鉛	鴛	塩	於	汚	甥	凹	央	奥	往	応					自

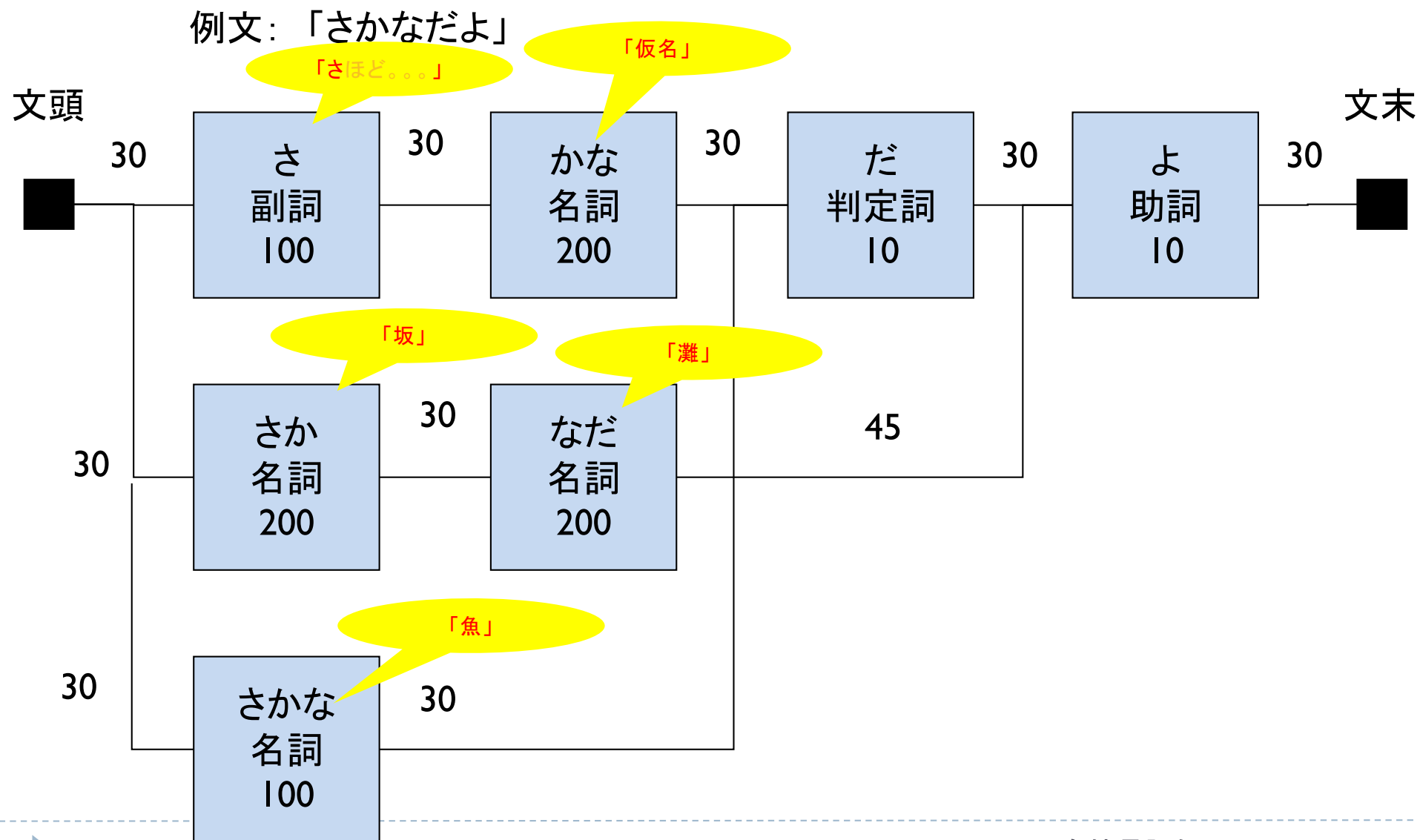
文字集合（Unicode）に対する文字符号化（UTF-8）

文字集合	文字符号化
U+0000~U+007F	00~7F (7bit)
(ASCIIと同一)	
U+0080~U+07FF	C280~DFBF (11bit)
(非漢字の一部)	
U+0800~U+FFFF	E08080~EFBFBF (16bit)
(残りの漢字)	
U+10000~U+1FFFFFF	F0808080~F7BFBFBF (21bit)
7 (第3、第4水準漢字)	自然言語処理 2023/5/4

形態素解析 解析結果の例

表層	読み	基本形	品詞	活用形	活用
いま	(いま)	いま	副詞		
米国	(べいこく)	米国	地名		
から	(から)	から	格助詞		
は	(は)	は	副助詞		
さまざまな	(さまざまな)	さまざま	形容詞	ナノ形容詞	ダ列基本連体
警報	(けいほう)	警報	普通名詞		
が	(が)	が	格助詞		
発せ	(はっせ)	発する	動詞	サ変動詞	文語未然形
られて	(られて)	られる	動詞性接尾辞	母音動詞	タ系連用テ形
いる	(いる)	いる	動詞性接尾辞	母音動詞	基本形
。	(。)	。	句点		
EOS					
商務	(しょうむ)	商務	普通名詞		
省	(しょう)	省	普通名詞		
の	(の)	の	接続助詞		
日本	(にっぽん)	日本	地名		
など	(など)	など	副助詞		
鉄鋼	(てっこう)	鉄鋼	普通名詞		
メーカー	(めーかー)	メーカー	普通名詞		
製品	(せいひん)	製品	普通名詞		
に	(に)	に	格助詞		
対する	(たいする)	対する	動詞	サ変動詞	基本形
ダンピング	(ダンピング)	ダンピング	カタカナ		
仮	(かり)	仮	普通名詞		
決定	(けってい)	決定	サ変名詞		
や	(や)	や	接続助詞		
上院	(じょういん)	上院	普通名詞		
.....					
EOS					

「茶釜」のコスト最小法



文脈自由文法の例 $G_k = (V_k, \Sigma_k, P_k, \text{文})$

$\Sigma_k = \{\text{きた, 文化, 伝わった, から, が}\}$

$V_k = \{\text{文, 後置詞句, 動詞, 名詞, 助詞}\}$
 $\cup \Sigma_k$

$P_k = \{ \text{文} \rightarrow \text{後置詞句 文},$
 $\text{文} \rightarrow \text{後置詞句 動詞},$
 $\text{後置詞句} \rightarrow \text{名詞 助詞},$
 $\text{後置詞句} \rightarrow \text{動詞 助詞},$
 $\text{名詞} \rightarrow \text{きた}, \text{名詞} \rightarrow \text{文化}$
 $\text{動詞} \rightarrow \text{きた}, \text{動詞} \rightarrow \text{伝わった}$
 $\text{助詞} \rightarrow \text{から}, \text{助詞} \rightarrow \text{が} \}$

導出の例

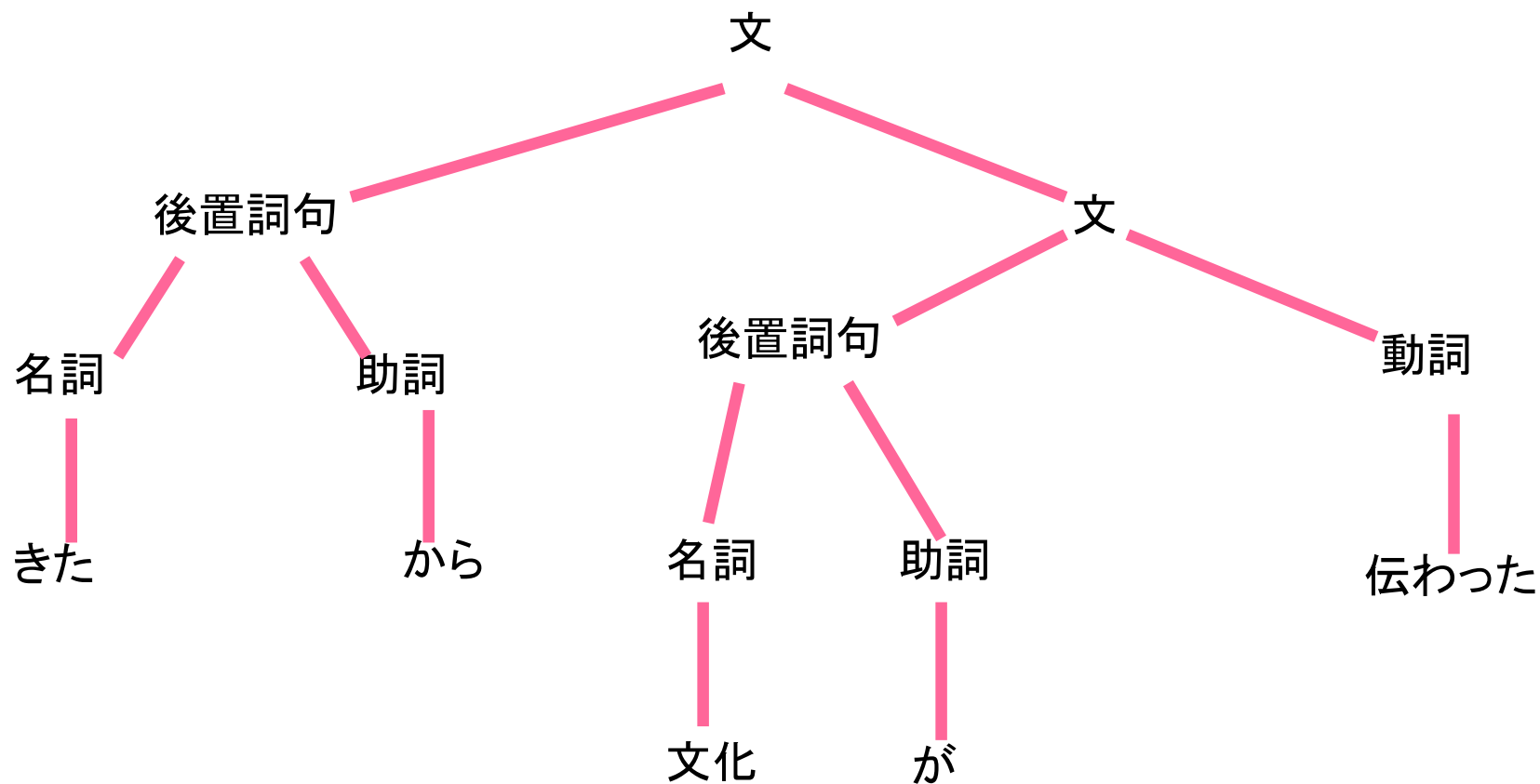
解析する文 w_k = きた から 文化 が 伝わった

導出

文 \Rightarrow 後置詞句 文
 \Rightarrow 名詞 助詞 文
 \Rightarrow きた 助詞 文
 \Rightarrow きた から 文
 \Rightarrow きた から 後置詞句 動詞
 \Rightarrow きた から 名詞 助詞 動詞
 \Rightarrow きた から 文化 助詞 動詞
 \Rightarrow きた から 文化 が 動詞
 \Rightarrow きた から 文化 が 伝わった

構文木 (例)

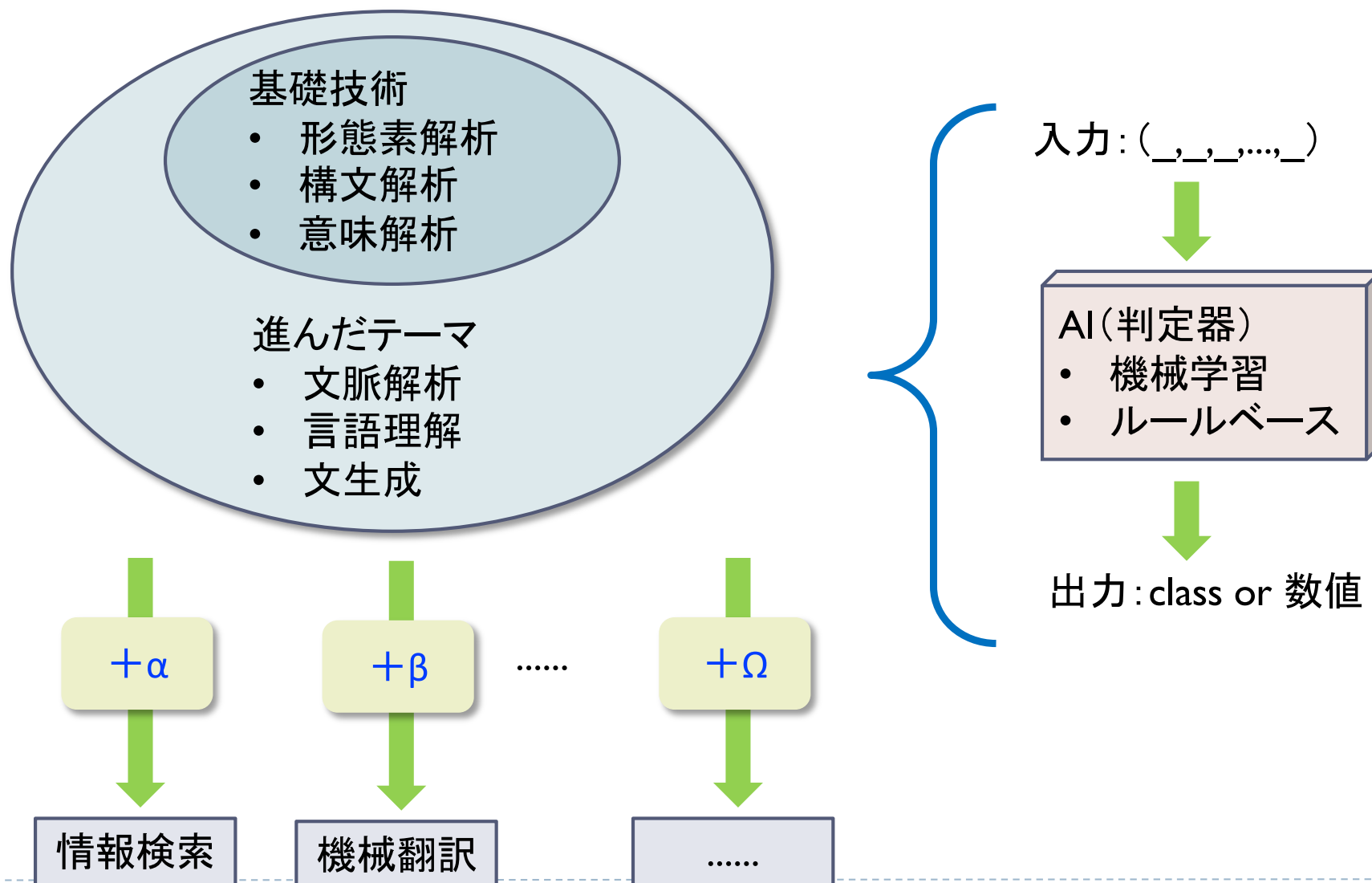
$w_k =$ きた から 文化 が 伝わった



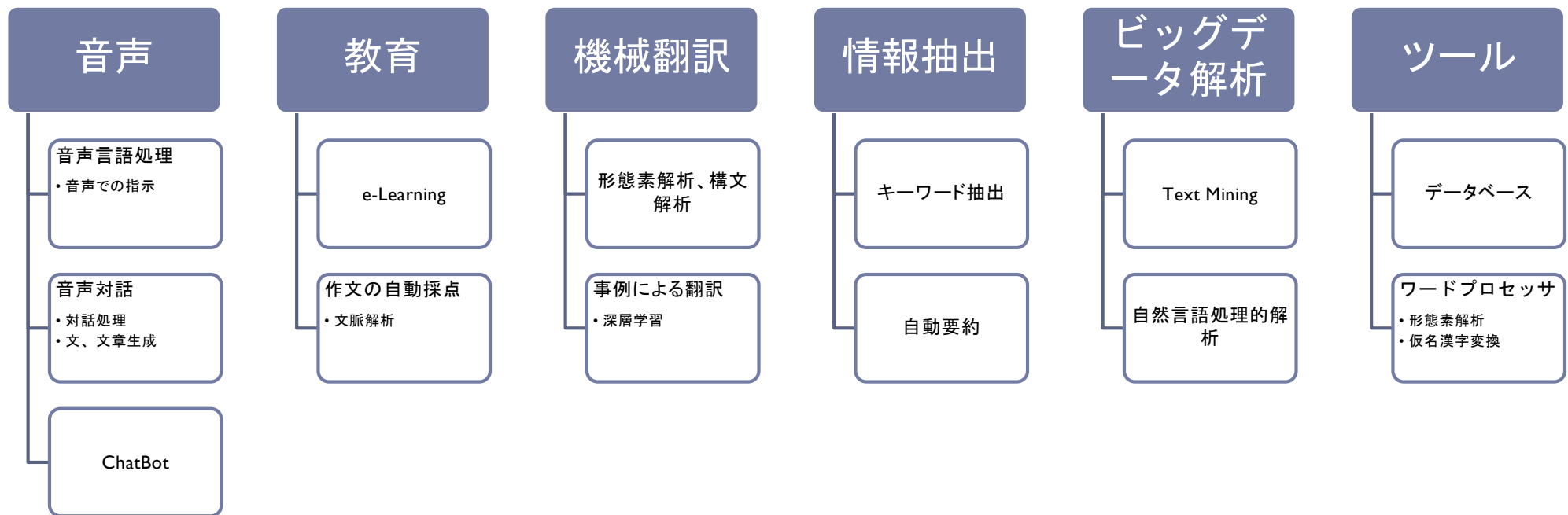
構文木(生成木、導出木)による文の構文構造の表現

「自然言語処理」に関連する各種技術

〜本講義で扱う各技術の位置づけ〜



自然言語処理の応用分野



1.2 自然言語処理関連の歴史

1940年代、W.Weaver, A.D.Booth

- ▶ 機械翻訳(英 \leftrightarrow 仏(France, 法))
- ▶ 単語の置き換えによる

1954年、ジョージタウン大、IBM

- ▶ ロシア語 \rightarrow 英語
- ▶ 直接翻訳方式。英語、ドイツ語、フランス語間である程度成功
 - step1) 単語レベルの置換。単語の活用形を全て展開
 - step2) 目的言語で語順入れ替え

1957年、N.Chomsky、“Syntactic Structure”

- ▶ 句構造文法
 - 文: 句構造(木構造)、変形により疑問文、受け身文などを文生成

1960年代

- ▶ 辞書、形態素解析、文法の導入 $\cdots \rightarrow$ 構文解析

予想分析法 (predictive analyzer)

- ▶ 文を途中まで聞いて構文構造を理解。予想しながら聞いている。

1961年、C.Hockett、Grammar for the Hearer

- ▶ Rhodes、ロシア語の解析文法
- ▶ 久野暲、英語文法

このような解析手法の研究が進んだが、構文的曖昧性が問題

They are flying planes.

Time flies like an arrow.

He saw a woman in the garden with a telescope.

日本

1955年～、九州大学、電気試験所

- ▶ 英独(徳)日間の機械翻訳

1959年、機械翻訳「やまと」(英→日)、九州大学(英独日)

- ▶ 漢字を扱えない→仮名
- ▶ 形態素解析が難しい -----> ワープロ(1978年)出現まで

1966年、米(美)、ALPAC報告書

- ▶ 機械翻訳は当分無理
- ▶ もっと基礎研究が必要

1960～1970前半、暗黒の時代

事務処理

1959年、IBM KWIC (H.P.Luhn)、検索結果の表示法

1964年、MEDLARS、医学文献の情報検索システム

1950年～1960年、オートマトン

意味解析の重要性

構文解析→精密

⇒ 構文的曖昧性→増加

⇒ 意味解析→重要

1968年、C.Fillmore、格文法(case grammar)

He opened the door by a key

A key opened the door

行為者(agent): he、対象(object): door、道具(instrument): key

それぞれの動詞について意味素(数十～数百個)

▶ 各動詞がどのような意味の名詞を格にとるか

⇒ 構文的曖昧性→小

情報検索

- ▶ 書誌的事項(本、論文の表題、著者名)→論文の抄録作成
- ▶ 1971年、MEDLINE(医学分野の世界最大のDB)

1976年～、日本科学技術情報センター

- ▶ 本文文書全体を検索語(書誌的事項、キーワード)でscan
- ▶ 全文検索(full text search)

人工知能(artificial intelligence)

1966年、ELIZA(対話システム)

⇒ 知識が重要

場面の分割が重要(T.Winograd)

代名詞の照応、曖昧さの解消、... 場面の認識により解消

計算言語学

1984年、M.Kay、functional unification grammar

- ▶ 隣接する2つの要素(単語)を一つ上の要素とする

LFG

- ▶ 2, 3の結合規則のみ + 単語に対して豊富な属性情報
-

言語理論の限界 ⇒ 大量の文章データ

1985年頃、英語フランス語の対訳を確率的に解析

→ 文法的、語彙的情報の抽出

1文単位の解析では不十分

- ▶ 文脈情報、場面知識

まず、隣接する文相互間の関係、解析

- ▶ 代名詞の照応、省略、話題、主題、焦点、
- ▶ 隣接する文間の関係(並列、理由、説明、...)

例文主導翻訳、実例型機械翻訳、用例翻訳

種々の特有の表現を例文(例句)として集め、規則化

AIの歴史

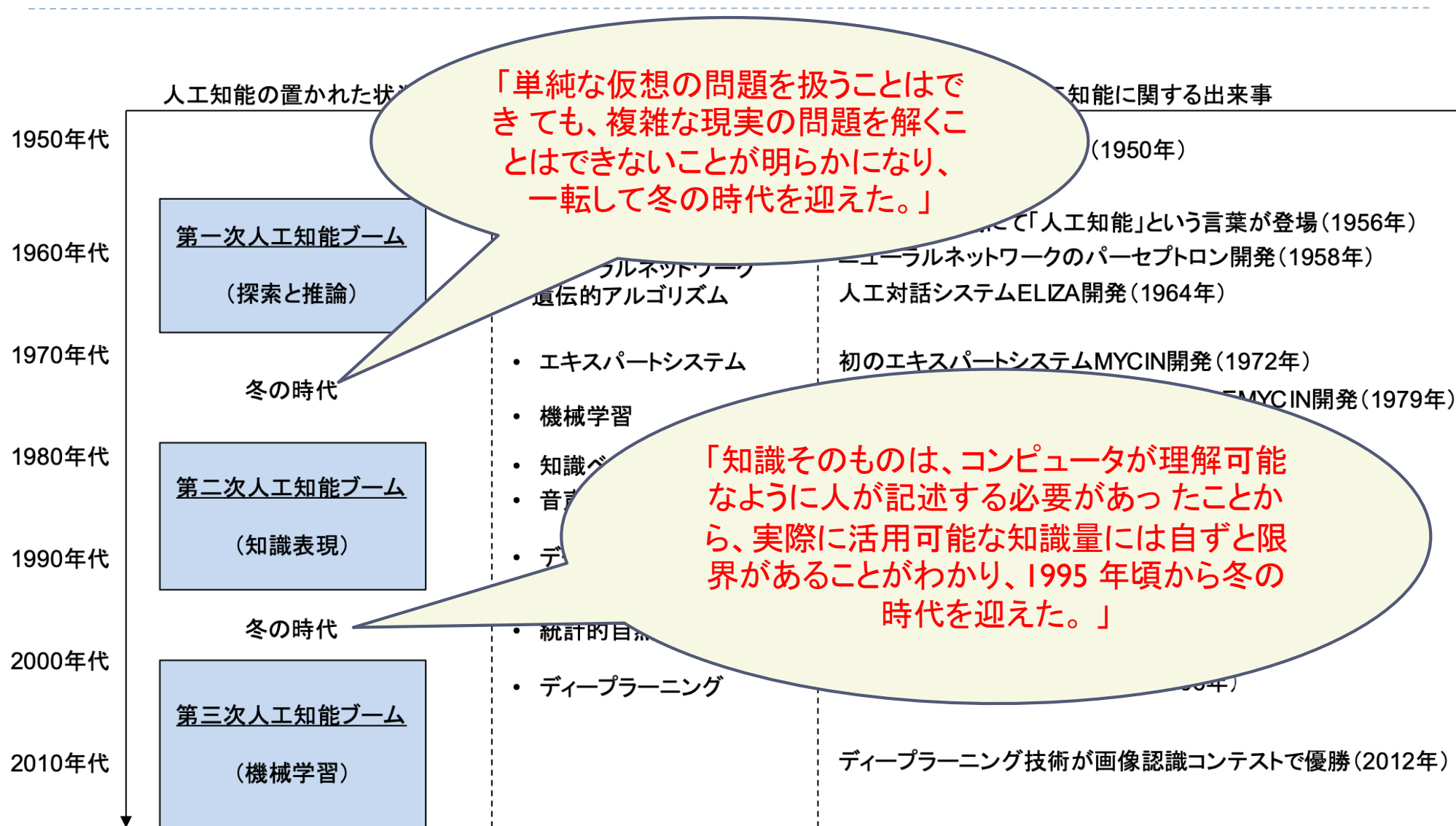
総務省「ITCの進化が雇傭と働き方に及ぼす調査研究(H28)

	人工知能の置かれた状況	主な技術等	人工知能に関する出来事
1950年代			チューリングテストの提唱(1950年)
1960年代	第一次人工知能ブーム (探索と推論)	<ul style="list-style-type: none">探索、推論自然言語処理ニューラルネットワーク遺伝的アルゴリズム	ダートマス会議にて「人工知能」という言葉が登場(1956年) ニューラルネットワークのパーセプトロン開発(1958年) 人工対話システムELIZA開発(1964年)
1970年代	冬の時代	<ul style="list-style-type: none">エキスパートシステム機械学習	初のエキスパートシステムMYCIN開発(1972年) MYCINの知識表現と推論を一般化したEMYCIN開発(1979年)
1980年代	第二次人工知能ブーム (知識表現)	<ul style="list-style-type: none">知識ベース音声認識	第五世代コンピュータプロジェクト(1982~92年) 知識記述のサイクプロジェクト開始(1984年) 誤差逆伝播法の発表(1986年)
1990年代	冬の時代	<ul style="list-style-type: none">データマイニングオントロジー	
2000年代	第三次人工知能ブーム (機械学習)	<ul style="list-style-type: none">統計的自然言語処理ディープラーニング	ディープラーニング技術の提唱(2006年)
2010年代			ディープラーニング技術が画像認識コンテストで優勝(2012年)



AIの歴史

総務省「ITCの進化が雇傭と働き方に及ぼす調査研究(H28)



歴史的な対話システム（AIの歴史）

- SHRDLU（1970年頃、テリー・ウィノグラード）
積み木の世界での積み木いじり操作に
限定することにより、人間とコン
ピュータとの会話システムを構築した。
映像（[SHRDLU in Action](https://www.youtube.com/watch?v=bo4RvYJYOzI), by Alec Julien
<https://www.youtube.com/watch?v=bo4RvYJYOzI>）
- ELIZA（1965年頃、MITのジョセフ・ワイゼンバウム）
DOCTORという来談者中心療法のセラピストを
シミュレート。
Emacs上で、M-x doctor
- MYCIN（1970年頃、スタンフォード大）
伝染性の血液疾患を診断し、抗生物質を
推奨するエキスパートシステム



まとめ

1. 自然言語処理の概要
 1. 自然言語処理の典型的モデル
 2. コンピュータ内での文字の扱い
 3. 形態素解析
 4. 構文解析
 5. 各技術の位置づけ
2. 自然言語処理の歴史

課題

算術式を記述できる文法 G_h により、次の式を構文木で表せ。

$$G_h = (\{S, T, F\} \cup \Sigma_h, \Sigma_h, P_h, S)$$

$$\Sigma_h = \{id, +, -, \times, \div, (,)\}$$

$$P_h = \{S \rightarrow S + T, S \rightarrow S - T, \\ S \rightarrow T, \\ T \rightarrow T \times F, T \rightarrow T \div F, \\ T \rightarrow F, \\ F \rightarrow id, \\ F \rightarrow (S)\}$$

1. $id + id + id$
2. $id \times id + id$
3. $id \times (id + id)$