

文法推論のコンパイラ言語への適用

Applying Grammatical Inference to a Compiled Language

常見 一樹 上野 敦志 田窪 朋仁 辰巳 昭治
Kazuki Tsunemi Atsushi Ueno Tomohito Takubo Shoji Tatsumi

大阪市立大学 大学院工学研究科
Osaka City University, Graduate School of Information Engineering

We aim to develop a method for mining grammatical rules of a programming language from a bunch of source codes and the corresponding byte codes. This method will lead to bridging the gap between the fields of grammatical inference and semantic understanding both of which are critical to language learning. We show the potential of our method by a thought experiment.

1. はじめに

プログラミング言語は、コンピュータに対する命令の記述をより分かりやすく行うために開発された人工言語の総称である。特に近年用いられているプログラミング言語は高級言語と呼ばれ、より人間の直感的に分かりやすい構文で記述できるように設計されている。ここで注目すべきは、人間にも事前知識さえあれば解読可能なプログラミング言語が、コンパイラやインタプリタ、ミドルウェア、OS、BIOSなどを経て最終的にCPUで処理することが出来る時系列シグナル情報に機械的に変換されるという点である。これはプログラミング言語という文法的に記述された構造情報が意味情報を保存したまま、手続き処理を纏めた構造へと変換されている事を意味する。しかもその変換は構文に関する情報さえ与えてやれば機械的に実行可能である。

構文に関する情報を用いてある構造を持つ表現から別の構造を持つ構造への変換が存在するとき(図1)、相異なる2つの構造で表現される同じ情報から変換に必要な情報を抜き出すことが可能であると考えられる。例えば、プログラミング言語のソースコードと、コンパイル後のオブジェクトコードからコンパイラ推定が可能である。

本論文ではソースコードとバイトコードを記号列のベクトルとして抽象化し、ソースコード同士及びバイトコード同士の差分をベクトル表現することで、推定に有用な情報を行列として抽出する手法を提案する。また構文の重要な機能の一つである構造の階層化に関する情報の抽出も試みる。

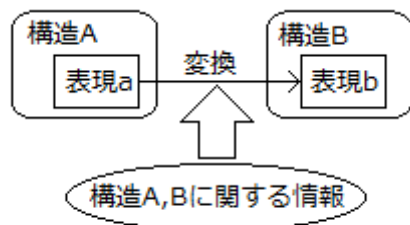


図 1

2. 関連研究

2.1 文法推論

法則性のある記号列を元にして、その記号列に共通して存在する表現拘束情報を抽出する技法は文法推論(Grammatical Inference)と呼ばれ盛んに研究されている。しかしこれは意味情報には殆ど言及せず、入力記号列が文法的に正か誤かという判断を行うのが主流である [Clark 2010]。

2.2 文脈自由文法

文脈自由文法は二型文法ともいい、その強力な表現能力と扱いやすさからコンパイラ言語の構文定義にも用いられている。自然言語文法の研究にも文脈自由文法を利用して文法の自動学習も試みられている [中村 2007]。また文脈自由文法中に出現する内部記号と概念に関する研究も行われている [花川 1998]。

2.3 自然言語処理

機械翻訳や文章からの知識情報の自動獲得は自然言語処理の分野で非常に盛んに研究されている。また応用も行われており、実際にサービスとして提供されている。しかしその品質は未だに非常に悪く、構文情報などを盛り込むことでの性能改善が期待されている [Nichols 2010]。既知の Wikipedia 独自のページ構造情報を利用する Wikipedia マイニングが一例である。 [玉川 2009]

3. 提案手法

Java ソースコードとそのソースコードをコンパイルした結果のバイトコードの対を1つの事例とし、この事例の集合を用いて構文に関する情報を抽出することを試みる。議論を簡単にするためにソースコードの字句解析が済んでいるものとし、バイトコードも1バイト単位でブロック化して扱う。

事例の有限集合を E とする。 i 番目の事例 e^i は字句配列ベクトル s^i とバイト配列ベクトル b^i を用いて以下のように表現される。ここで T は転置を表す。 $|E|$ は集合 E の要素数を表す。

$$e^i = (s^i, b^i)^T, \quad e^i \in E \quad (i = 1, 2, \dots, |E|) \quad (1)$$

また字句配列ベクトル s^i とバイト配列ベクトル b^i はそれぞれ、ソースコードを成す字句の集合 C とバイトデータの集合 N を用いて以下のように表される。

連絡先: 常見一樹, 大阪市立大学大学院工学研究科,
〒558-8585 大阪市住吉区杉本 3-3-138
TEL: 06-6605-2688
Mail: ksilver@kdel.info.eng.osaka-cu.ac.jp

$$\mathbf{s}^i = (s_1^i, s_2^i, \dots, s_{l_i}^i)^T, \quad \mathbf{s}^j \in \mathbf{C} \ (j = 1, 2, \dots, l) \quad (2)$$

$$\mathbf{b}^i = (b_1^i, b_2^i, \dots, b_{m_i}^i)^T, \quad \mathbf{b}^j \in \mathbf{N} \ (j = 1, 2, \dots, m) \quad (3)$$

ここで l_i と m_i はそれぞれ \mathbf{s}^i と \mathbf{b}^i の次元の大きさを表す。

字句配列ベクトルの集合を $\mathbf{S} = \{\mathbf{s}^i | \mathbf{e}^i \in \mathbf{E}\}$, バイト配列ベクトルの集合を $\mathbf{B} = \{\mathbf{b}^i | \mathbf{e}^i \in \mathbf{E}\}$ とする。コンパイラ推定問題は、コンパイラによる変換を f とすると、 $f(\mathbf{s}^i) = \mathbf{b}^i \ (i = 1, 2, \dots, |\mathbf{E}|)$ を満たす変換 $f: \mathbf{S} \rightarrow \mathbf{B}$ を求める問題と言うことが出来る。

3.1 構造の作成

(1) 次元一致ケース

相異なる2つの事例 $\mathbf{e}^p, \mathbf{e}^q$ について、字句配列ベクトル $\mathbf{s}^p, \mathbf{s}^q$ の次元とバイト配列ベクトル $\mathbf{b}^p, \mathbf{b}^q$ の次元について以下の式が成り立つ場合を考える。

$$l_p = l_q, \quad m_p = m_q$$

相異なる2つの事例 $\mathbf{e}^p, \mathbf{e}^q$ に対して事例差 $\Delta \mathbf{e}^{p,q}$, 及び字句差異ベクトル $\Delta \mathbf{s}^{p,q}$ とバイト差異ベクトル $\Delta \mathbf{b}^{p,q}$ を以下のように定義する。

$$\Delta \mathbf{e}^{p,q} = (\Delta \mathbf{s}^{p,q}, \Delta \mathbf{b}^{p,q})^T \quad (4)$$

$$\Delta \mathbf{s}^{p,q} = (\Delta s_1^{p,q}, \Delta s_2^{p,q}, \dots, \Delta s_{l_p}^{p,q})^T$$

$$\Delta s_i^{p,q} = \begin{cases} 1 & (s_i^p \neq s_i^q) \\ 0 & (s_i^p = s_i^q) \end{cases}, \quad (i = 1, 2, \dots, l_p) \quad (5)$$

$$\Delta \mathbf{b}^{p,q} = (\Delta b_1^{p,q}, \Delta b_2^{p,q}, \dots, \Delta b_{m_p}^{p,q})^T$$

$$\Delta b_i^{p,q} = \begin{cases} 1 & (s_i^p \neq s_i^q) \\ 0 & (s_i^p = s_i^q) \end{cases}, \quad (i = 1, 2, \dots, m_p) \quad (6)$$

ここで $\Delta \mathbf{s}^{p,q}$ と $\Delta \mathbf{b}^{p,q}$ はそれぞれソースコード情報とバイトデータ情報の差分を表す。

$\Delta \mathbf{s}^{p,q}$ が零ベクトルの時、 $\mathbf{s}^p = \mathbf{s}^q$ である。コンパイラは決定性を持つため $\mathbf{b}^p = \mathbf{b}^q$ である。このとき $\mathbf{e}^p = \mathbf{e}^q$ となり $\mathbf{e}^p, \mathbf{e}^q$ が相異なる事例であることに反する。よって $\Delta \mathbf{s}^{p,q}$ は零ベクトルになり得ない。 $\Delta \mathbf{b}^{p,q}$ が零ベクトルの時、コンパイラによる変換 f について、 $\mathbf{b}^p = \mathbf{b}^q$ より $f(\mathbf{s}^p) = f(\mathbf{s}^q)$ が言える。

$\Delta \mathbf{b}^{p,q}$ が零ベクトルでない場合、 $\Delta \mathbf{e}^{p,q}$ はソースコードの違いによるバイトデータ情報の変化が記録されていることになる。ここで長さが 1 になるように $\Delta \mathbf{s}^{p,q}, \Delta \mathbf{b}^{p,q}$ をそれぞれ正規化した正規字句差異ベクトル $\delta \mathbf{s}^{p,q}$ と、正規バイト差異ベクトル $\delta \mathbf{b}^{p,q}$ を定義する。

$$\delta \mathbf{s}^{p,q} = \frac{\Delta \mathbf{s}^{p,q}}{\|\Delta \mathbf{s}^{p,q}\|}, \quad \delta \mathbf{s}^{p,q} = (\delta s_1^{p,q}, \delta s_2^{p,q}, \dots, \delta s_{l_p}^{p,q})^T \quad (7)$$

$$\delta \mathbf{b}^{p,q} = \frac{\Delta \mathbf{b}^{p,q}}{\|\Delta \mathbf{b}^{p,q}\|}, \quad \delta \mathbf{b}^{p,q} = (\delta b_1^{p,q}, \delta b_2^{p,q}, \dots, \delta b_{m_p}^{p,q})^T \quad (8)$$

$\delta \mathbf{s}^{p,q}, \delta \mathbf{b}^{p,q}$ はそれぞれ差異が複数箇所ある場合に、差異の箇所が増えるに従い各成分の値が小さくなる性質を持つ、全体に対する差分量を表す量である。ソースコードの差異情報とバイトデータの差異情報の対応関係を表す構造対応行列 $\mathbf{M}^{p,q}$ を以下のように定義する。

$$\mathbf{M}^{p,q} = \delta \mathbf{s}^{p,q} (\delta \mathbf{b}^{p,q})^T \quad (9)$$

ここで \mathbf{s}^p と \mathbf{s}^q , \mathbf{b}^p と \mathbf{b}^q が共に似ている場合、 $\delta \mathbf{s}^{p,q}, \delta \mathbf{b}^{p,q}$ は多くの成分が 0 で少数成分が 1 に近いベクトルとなる。このとき行列 $\mathbf{M}^{p,q}$ も定義から同様に少数成分が 1 に近い構造となる。この

とき $\mathbf{M}^{p,q}$ の (i, j) 成分が非零であるとき、 $\delta \mathbf{s}^{p,q}$ の第 i 成分と $\delta \mathbf{b}^{p,q}$ の第 j 成分が共に非零である。これは事例 $\mathbf{e}^p, \mathbf{e}^q$ において、 \mathbf{s}^p の第 i 番目と \mathbf{b}^p の第 j 成分、 \mathbf{s}^q の第 i 番目と \mathbf{b}^q の第 j 成分が対応していることを示す。また (i, j) 成分の大きさが対応関係の強さを表している。この値を対応度と呼ぶことにする。対応度は \mathbf{s}^p と \mathbf{s}^q , \mathbf{b}^p と \mathbf{b}^q が似ていない場合でも求めることが出来るが値は小さくなる。つまり $\mathbf{M}^{p,q}$ は $\mathbf{e}^p, \mathbf{e}^q$ から求まる対応関係の位置と関連度を表す行列である。

多くの事例ペアから同次元の構造対応行列が得られ、なおかつ特定の成分 (i, j) が共通して 1 に近い時、ソースコードの第 i 番目の字句とバイトコードの第 i 番目のバイトが対応している可能性が高い。この字句とバイトの対応はコンパイラを構成する上で非常に重要である。

(2) 次元不一致ケース

相異なる2つの事例 $\mathbf{e}^p, \mathbf{e}^q$ について、字句配列ベクトル $\mathbf{s}^p, \mathbf{s}^q$ の次元が異なる。もしくはバイト配列ベクトル $\mathbf{b}^p, \mathbf{b}^q$ の次元が異なる場合を考える。

字句配列ベクトル $\mathbf{s}^p, \mathbf{s}^q$ について、 $\dim(\mathbf{s}^p) > \dim(\mathbf{s}^q)$ であるとする。ただし $\dim(\mathbf{s})$ はベクトルの次元を表す記号である。

$$\mathbf{s}^p = (s_1^p, s_2^p, \dots, s_{l_p}^p)^T, \quad s_i^p \in \mathbf{C}, \quad i = 1, 2, \dots, l_p \quad (10)$$

$$\mathbf{s}^q = (s_1^q, s_2^q, \dots, s_{l_q}^q)^T, \quad s_i^q \in \mathbf{C}, \quad i = 1, 2, \dots, l_q \quad (11)$$

$$l_p > l_q \quad (12)$$

このとき字句配列圧縮記号 X を用いて第 k 成分から $\dim(\mathbf{s}^p) - \dim(\mathbf{s}^q) + 1$ 個を置き換えた抽象字句配列ベクトル $\mathbf{s}^{p'}$ を作る。

$$\mathbf{s}^{p'} = (s_1^p, s_2^p, \dots, s_{k-1}^p, X, s_{k+n-m+1}^p, \dots, s_{l_p}^p)^T \quad (13)$$

$$X = (s_k^p, s_{k+1}^p, \dots, s_{k+n-m}^p)^T \quad (14)$$

$$\dim(\mathbf{s}^{p'}) = \dim(\mathbf{s}^q) \quad (15)$$

字句圧縮記号 X は一意に定義できないため、 $\mathbf{s}^{p'}, \mathbf{s}^q$ から得られる字句差異ベクトル $\Delta \mathbf{s}^{p,q}$ の長さが最小になるように選ぶ。バイト配列ベクトル $\mathbf{b}^{p'}, \mathbf{b}^q$ に対しても同様に、バイト配列圧縮記号を用いて抽象バイト配列ベクトルを作る。すると $\dim(\mathbf{s}^{p'}) = \dim(\mathbf{s}^q)$ かつ $\dim(\mathbf{b}^{p'}) = \dim(\mathbf{b}^q)$ であるため $\mathbf{e}^{p'} = (\mathbf{s}^{p'}, \mathbf{b}^{p'})$, $\mathbf{e}^q = (\mathbf{s}^q, \mathbf{b}^q)$ について、次元一致ケースの手続きを行うことが出来る。

3.2 統合と一般化

前節の手続きだけでは抽出した情報は事例のペアに依存しているため、事例のペア毎に構文対応行列や字句圧縮記号といった構造情報が生み出されることになる。事例の数に対して組み合わせの爆発が生じ、分析した情報の急激な増加することが危惧される。

そこで多数生み出される構造情報のうち比較的似ている構造情報や字句圧縮記号、バイト圧縮記号を等価と見なす。複数の構造情報から共通する性質を抽出し、構造情報を一般化するなどの操作を加える。こうすることで分析情報の圧縮と表現の一般化を同時に行う事が出来る。

構造対応行列が直交行列となっており、記号の対応度が 1 で有れば、等価変換と行列演算だけでコードの書き換えが可能である。例えば未知のソースコード $\mathbf{s} = (x, *, y)^T$ を構造対応行列 \mathbf{M} で書き換えることを考える。

$$Ms = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} x \\ * \\ y \end{pmatrix} = \begin{pmatrix} x \\ y \\ * \end{pmatrix} \quad (16)$$

これはポーランド記法の変換である。実際の Java バイトコードにおいても変数名格納場所などは別であるものの、load_1 命令→load_2 命令→imul 命令など似たような構造をとる。多くの事例ペアにおける構造対応行列に対し平均を取るなどの操作を行うことで、ソースコード中の位置とバイトコードの位置の対応関係が予測できる。つまり構文構造の対応関係の発見に繋がる。

また圧縮記号の導入も単に次元の調整ではなく、一般化プロセスでは、文法の子構造を検出するのに役立つと考えられる。例えば字句圧縮記号は字句配列ベクトルと同じ構造を取る。他の字句記号配列ベクトルと字句圧縮記号の間に等価関係や構造情報が抽出できた際に字句記号配列ベクトルは子構造とみなす事が出来る。以下の例は if 文的構造から入れ子にされた式の構造を抽出するプロセスである。ソースコード s^p, s^q を以下のように定義する。

$$s^p = (if, a, \neq, b, then)^T, \quad s^q = (if, bool, then)^T \quad (17)$$

ここで字句圧縮記号 $A = (a, \neq, b)^T$ を置くと、

$$s^{p'} = (if, A, then)^T, \quad s^{q'} = (if, bool, then)^T \quad (18)$$

ここから A と $bool$ 、即ち $a \neq b$ と $bool$ が同じ種類のクラスに分類されることが推察できる。

4. 実験

Java を含め近年開発されたプログラミング言語は、人間が読みやすく、なおかつ汎用的である。その代償としてコンパイラの仕様は複雑になっており、これは本研究で扱うコンパイル時の構文情報もまた煩雑であることを意味する。推定対象の構文が複雑になればなるほど、実験するのに必要な計算量とデータセットが爆発的に増大するのは明白である。同時に、プログラミング言語全体という広域の推定結果を実験的に求めることが出来たととしても、内部でどのような処理が行われたか分析は困難である。

そこで本論文では限定された部位の構文情報の推定を行い、提案手法の有用性及び、その挙動を分析する。具体的な初期段階としては、Java コードのうちクラス定義を固定化し、特定メソッド内部の手続きのみを変更して構文情報の推定を行う。

現在は準備段階として、式の間演算子表現とポーランド記法間の変換規則の推定を行っている。これはメソッドの手続き記述の中でも最も初歩的な機能の一つである式の計算において、Java のバイトデータはポーランド記法に非常に類似した構造を出力するためである。中間演算子表現からポーランド記法への変換に関する情報の抽出例を以下に述べる。

中間演算子表現からポーランド記法表現の変換について考える。以下に挙げる事例に対して提案手法を用いた場合に得られる構造情報について具体的に述べる。中間演算子表現をソースコードとして扱い、ポーランド記法をバイトコードとして扱う。

4.1 次元一致ケース

(事例 1) ソースコード: $x + y$, バイトコード: (X, Y, Add)

(事例 2) ソースコード: $x - y$, バイトコード: $(X, Y, Minus)$

(事例 3) ソースコード: $z - y$, バイトコード: $(Z, Y, Minus)$

記号 e^1, e^2, e^3 を用いて表現すると以下ようになる。

$$e^1 = (s^1, b^1)^T, \quad s^1 = (x, +, y)^T, \quad b^1 = (X, Y, Add)^T \quad (19)$$

$$e^2 = (s^2, b^2)^T, \quad s^2 = (x, -, y)^T, \quad b^2 = (X, Y, Minus)^T \quad (20)$$

$$e^3 = (s^3, b^3)^T, \quad s^3 = (z, -, y)^T, \quad b^3 = (Z, Y, Minus)^T \quad (21)$$

e^1, e^2 を用いて差分をとる。

$$\Delta s^{1,2} = (0, 1, 0)^T, \quad \Delta b^{1,2} = (0, 0, 1)^T \quad (22)$$

正規化する。

$$\delta s^{1,2} = (0, 1, 0)^T, \quad \delta b^{1,2} = (0, 0, 1)^T \quad (23)$$

構造対応行列を求める

$$M^{1,2} = \delta s^{1,2} (\delta b^{1,2})^T = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} \quad (24)$$

(2,3)成分が 1 だから s^1, s^2 第 2 成分と b^1, b^2 第 3 成分を比較して、字句 "+" とコード "Add", 字句 "-" とコード "Minus" がそれぞれ等しいことが分かる。同様に事例 e^2, e^3 については以下の構造対応行列が得られる。

$$M^{2,3} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (25)$$

(1,1)成分が 1 だから s^2, s^3 第 1 成分と b^2, b^3 第 1 成分を比較して、字句 "x" とコード "X", 字句 "y" とコード "Y" がそれぞれ等しいことが分かる。同様に事例 e^3, e^1 については正規字句配列ベクトルと正規バイト配列ベクトルが以下ようになる。

$$\delta s^{3,1} = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}, 0\right)^T, \quad \delta b^{3,1} = \left(\frac{1}{\sqrt{2}}, 0, \frac{1}{\sqrt{2}}\right)^T \quad (26)$$

このため構造変換行列は以下の通りになる。

$$M^{3,1} = \begin{pmatrix} 1/2 & 0 & 1/2 \\ 1/2 & 0 & 1/2 \\ 0 & 0 & 0 \end{pmatrix} \quad (27)$$

この構造変換行列には 1 の成分がないため対応関係は決定できない。例えば (1,1)成分に注目すると対応度が 0.5 で、対応する成分が "x" とコード "X", 字句 "y" とコード "Y" の対応を示している。これは正しい対応である。一方で (1,3)成分に注目すると同じく対応度が 0.5 で、対応する成分が "x" とコード "Add", 字句 "z" とコード "Minus" の対応を示している。これは誤った対応である。

このように対応度が 1 の場合は正しい対応が得られる一方で対応度が 1 より小さい場合は必ずしも正しい対応のみを表すとは限らない。

4.2 次元不一致ケース

(事例 1) ソースコード: $x + y$, バイトコード: (X, Y, Add)

(事例 2) ソースコード: $x + y * z$,

バイトコード: (X, Y, Z, Mul, Add)

記号 e^1, e^2 を用いて表現すると以下ようになる。

$$e^1 = (s^1, b^1)^T, \quad s^1 = (x, +, y)^T, \quad b^1 = (X, Y, Add)^T \quad (28)$$

$$e^2 = (s^2, b^2)^T$$

$$s^2 = (x, +, y, *, z)^T, \quad b^2 = (X, Y, Z, Mul, Add)^T \quad (29)$$

$\|\Delta s^{1,2'}\|, \|\Delta b^{1,2'}\|$ がそれぞれ最小になるように字句圧縮記号 a , バイト圧縮記号 A , 及びを用いて事例 e^2 を書き換えたに事例 $e^{2'}$ を作る。

$$e^{2'} = (s^{2'}, b^{2'})^T, \quad s^{2'} = (x, +, a)^T, \quad b^{2'} = (X, A, Add)^T \quad (30)$$

$$a = (y, *, z)^T, \quad A = (Y, Z, Mul)^T \quad (31)$$

次元一致ケースと同様に計算する。

$$\delta s^{1,2'} = (0, 0, 1)^T, \quad \delta b^{1,2'} = (0, 1, 0)^T \quad (32)$$

$$M^{1,2'} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \quad (33)$$

(3,2)成分が1だから s^1, s^2 第3成分と b^1, b^2 第2成分を比較して, 字句"y"とコード"Y", 字句"y,*,z"とコード"Y,Z,Mul"がそれぞれ等しいことが分かる. 字句"y,*,z"とコード"Y,Z,Mul"のペアは数式で言うと式の子構造である項と見なすことが出来る.

5. おわりに

今回提案した手法ではソースコードとバイトコードをベクトル化することで差分情報をベクトルとして求め, コンパイル前後の位置の対応関係を行列で表現した. その行列を用いて特別な事例ペアでは行列から一対一の対応が得られること, 任意のペアから対応関係を定量的に導出できることを示した.

また比較するソースコードのベクトルとしての次元が一致しない, バイトコードのベクトルとしての次元が一致しない場合は圧縮記号を用いて任意の文字列を一纏めにして長さを縮める手法を導入し, スケールの変化に対して柔軟に対応できる可能性も示した. これは構文の階層構造を探る上で重要である.

抽出した構造情報を組み合わせるなどして適用範囲を広げるなどの一般化に関する操作をどのように行うかが今後の課題である. また実際に提案手法を実装し, どのような規模でどのような変換データが得られるのかの実験的確認が必要である.

参考文献

- [Clark 2010] Alexander Clark, Christophe Costa Florêncio, Chris Watkins: "Languages as hyperplanes: grammatical inference with string kernels", Mach Learn, 2011.
- [中村 2007] 中村 克彦, 杉田 雄大: "文脈自由文法の漸次学習方式とその応用", 東京電機大学総合研究所年報, 2007.
- [花川 1998] 花川 賢治, 武田 英明, 西田 豊明: "文章から概念を獲得するための文法推論", 電子情報通信学会, 1998.
- [Nichols 2010] Nichols: "Applying Deep Grammars to Machine Translation, Paraphrasing, and Ontology Construction", NAIST-IS-DD0561041, NAIST, 2010.
- [玉川 2009] 玉川 奨, 桜井 慎弥, 手島 拓也: "日本語 Wikipedia からの大規模オントロジー学習", 人工知能学会誌 25 巻 5 号 SP-C, 2010.