

四書を学んだ MeCab + UDPipe は センター試験の漢文を読めるのか

安岡孝一*

1 はじめに

筆者が班長を務める京都大学人文科学研究所共同研究班「東アジア古典文献コーパスの実証研究」(班員: ウィッテルン クリストイアン、守岡知彦、池田巧、山崎直樹、二階堂善弘、鈴木慎吾、師茂樹、李媛、白須裕之、藤田一乗)では、現在、古典中国語(漢文)の依存文法解析に精力を傾注しており、その道具立ての一つとして、Universal Dependencies(以下「UD」)[1]の古典中国語への適用を研究している。依存文法解析それ自体は、Tesnière の構造的統語論[2]に源を発し、Мельчук の有向グラフ記述[3]によって、一応の完成を見た手法である。その最大の特長は、言語横断的な記述が可能だという点にあり、Мельчук の手法をコンピュータ向けに洗練した UDにおいても、言語に関わらない記述、という特長が前面に押し出されている。UDにおける文法構造記述は、句構造を考慮せず、全てを単語間のリンクとして表現する。これは、Мельчук の有向グラフ記述が、単語間のリンクという形態を取っていたからであり、そういう割り切りの結果として、言語横断的な文法構造記述を可能としているのである。

しかしながら、古典中国語の白文は、単語の間にも文の間にも区切りを持たず、その点で、欧米諸語とは異なる処理手法を必要とする。すなわち、白文の解析においては、まず、単語の区切りを認識することが必須であり、そのためには形態素解析をおこなわねばならない。入力された白文に対し、形態素解析[4-13]により、単語切りをおこなうと同時に、各単語の品詞を得る。その後に、依存文法解析[14-19]により、係り受け関係(単語間のリンク)を解析すると同時に、文の切れ目を得る、という手順になる。

これまでにわれわれは、古典中国語の形態素解析に MeCab [20] を用いてきた。MeCab は、もともとは日本語向けの形態素解析エンジンだったが、言語、辞書、コーパスに依存しない汎用的な設計がなされており、辞書とコーパスを準備すればいける言語にも対応できる。一方、古典中国語の依存文法解析には、UDPipe [21] が向いているようである。UDPipe は、チェコ語の係り受け解析エンジンが土台となっているが、UDへの設計拡張をおこなった際に、単語切り・品詞付与・文切りの機能が追加され、形態素解析と依存文法解析の両方を兼ね備えたものとなっている。ただ、UDPipe の形態素解析機能が、MeCab の形態素解析性能を凌駕しているのかどうかは、多少、疑問がある。

この問題に対し、われわれは、MeCab の形態素解析機能を、どのように UDPipe の依存文法解析機能と接続していくかについて、いくつかの比較実験をおこなった。学習用のコーパスは四書(孟子・論語・大学・中庸)を用いた。実験対象は、この 5 年間に実施された大学入試センター試験『国語』の本試験から、第 4 問の本文を用いた。また、比較実験の結果に対して、多少の改善を試みた。以下、実験と改善の概要を、順に述べる。

*京都大学人文科学研究所附属東アジア人文情報学研究センター

2 MeCab + UDPipe によるセンター漢文の解析

白文の形態素解析および依存文法解析は、単語切り・品詞付与・係り受け解析・文切りの4段階から成る。MeCab 0.996と UDPipe 1.2.0を、これらの4段階に適用するにあたって、われわれは、以下の3つの手法を比較することにした。

手法① 単語切り+品詞付与+係り受け解析+文切りの全てに UDPipe を用いる

手法② 単語切りに MeCab を、品詞付与+係り受け解析+文切りに UDPipe を用いる

手法③ 単語切り+品詞付与に MeCab を、係り受け解析+文切りに UDPipe を用いる

MeCab と UDPipe は、いずれも、Kanripo [22] をもとに手作業で作成した孟子・論語・大學・中庸の全文 UD コーパス(合計 56767 字、異なり字数 2333 字、文字コードは UTF-8)を学習済みである。全文 UD コーパスの一部を、以下に示す。

```
# newdoc id = KR1h0001_001
# text = 梁惠王上
1 梁 梁 PROPN n,名詞,主体,国名 Case=Loc|NameType=Nat 3 nmod   - Gloss=[country-name]|SpaceAfter=No
2 惠 惠 PROPN n,名詞,人,その他の人名 NameType=Prs      3 compound - Gloss=Hui|SpaceAfter=No
3 王 王 NOUN  n,名詞,人,役割          0 root    - Gloss=king|SpaceAfter=No
4 上 上 NOUN  n,名詞,固定物,関係       - Case=Loc      3 list    - Gloss=up|SpacesAfter=\n

# newpar
# text = 孟子見梁惠王
1 孟子 孟子 PROPN n,名詞,人,複合の人名 NameType=Prs      2 nsubj   - Gloss=Mencius|SpaceAfter=No
2 見 見 VERB v,動詞,行為,動作          0 root    - Gloss=see|SpaceAfter=No
3 梁 梁 PROPN n,名詞,主体,国名 Case=Loc|NameType=Nat 5 nmod   - Gloss=[country-name]|SpaceAfter=No
4 惠 惠 PROPN n,名詞,人,その他の人名 NameType=Prs      5 compound - Gloss=Hui|SpaceAfter=No
5 王 王 NOUN  n,名詞,人,役割       -           2 obj    - Gloss=king|SpaceAfter=No

# text = 王曰
1 王 王 NOUN  n,名詞,人,役割          -           2 nsubj   - Gloss=king|SpaceAfter=No
2 曰 曰 VERB v,動詞,行為,伝達        -           0 root    - Gloss=say|SpaceAfter=No
```

MeCab 辞書には、これら四書に出現する単語(形態素)に加え、図1左のように、固有名詞などを大量に追加[11]している。さらに、この MeCab 辞書を変形し、UDPipe の外部辞書(図1右)にも用いることで、できるだけ実験条件を揃えるようにした。UDPipe の文切りは、`--tokenizer=joint_with_parsing`オプションにより、係り受け解析と連動させている。上記の3手法の比較実験対象としては、過去5年間のセンター試験『国語』

| | |
|--|---|
| 江,0,0,0,n,名詞,人,名,*,*,江,*,*,[given-name] | 江 江 PROPN n,名詞,人,名 NameType=Giv |
| 江,0,0,0,n,名詞,人,姓氏,*,江,*,*,[surname] | 江 江 PROPN n,名詞,人,姓氏 NameType=Sur |
| 江,0,0,0,n,名詞,固定物,地名,*,*,江,*,*,[place-name] | 江 江 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江中,0,0,0,n,名詞,固定物,地名,*,*,江中,*,*,[place-name] | 江中 江中 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江北,0,0,0,n,名詞,固定物,地名,*,*,江北,*,*,[place-name] | 江北 江北 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江南,0,0,0,n,名詞,固定物,地名,*,*,江南,*,*,[place-name] | 江南 江南 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江壅,0,0,0,n,名詞,固定物,地名,*,*,江壅,*,*,[place-name] | 江壅 江壅 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江夏,0,0,0,n,名詞,固定物,地名,*,*,江夏,*,*,[place-name] | 江夏 江夏 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江寧,0,0,0,n,名詞,固定物,地名,*,*,江寧,*,*,[place-name] | 江寧 江寧 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江州,0,0,0,n,名詞,固定物,地名,*,*,江州,*,*,[place-name] | 江州 江州 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江戸,0,0,0,n,名詞,固定物,地名,*,*,江戸,*,*,[place-name] | 江戸 江戸 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江戸,0,0,0,n,名詞,固定物,地名,*,*,江戸,*,*,[place-name] | 江戸 江戸 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江東,0,0,0,n,名詞,固定物,地名,*,*,江東,*,*,[place-name] | 江東 江東 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江浙,0,0,0,n,名詞,固定物,地名,*,*,江浙,*,*,[place-name] | 江浙 江浙 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江淮,0,0,0,n,名詞,固定物,地名,*,*,江淮,*,*,[place-name] | 江淮 江淮 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江西,0,0,0,n,名詞,固定物,地名,*,*,江西,*,*,[place-name] | 江西 江西 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江進,0,0,0,n,名詞,人,名,*,*,江進,*,*,[given-name] | 江進 江進 PROPN n,名詞,人,名 NameType=Giv |
| 江陵,0,0,0,n,名詞,固定物,地名,*,*,江陵,*,*,[place-name] | 江陵 江陵 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江都,0,0,0,n,名詞,固定物,地名,*,*,江都,*,*,[place-name] | 江都 江都 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |
| 江都,0,0,0,n,名詞,固定物,地名,*,*,江都,*,*,[place-name] | 江都 江都 PROPN n,名詞,固定物,地名 Case=Loc NameType=Geo |

図 1: MeCab 辞書(左)と UDPipe の外部辞書(右)における固有名詞の例

の漢文(図3・23・43・63・83)から、句読点や返り点などを除去し、完全な白文を準備した。なお、センター試験の漢文は常用漢字で書かれているが、白文の字体を変更するようなことはせず、その代わり、MeCab 辞書と UDPipe 外部辞書に常用漢字の単語を追加した。処理結果を図4～6・24～26・44～46・64～66・84～86 に示す。

これらの処理結果を定量的に比較すべく、LAS (Labeled Attachment Score)・MLAS (Morphology-aware Labeled Attachment Score)・BLEX (Bi-LEXical dependency score) の3つの指標[23]を用いることにした。ただし、これらの指標を用いるためには、問題に対する正しい答が必要となることから、手作業で「正解」UD(図22・42・62・82・102)を準備した。比較結果を表1に示す。

表1: 手法①②③の LAS/MLAS/BLEX

| | 2019年 | 2018年 | 2017年 | 2016年 | 2015年 |
|-----|-------------------|-------------------|-------------------|-------------------|-------------------|
| 手法① | 45.48/39.87/43.09 | 42.86/37.54/38.77 | 49.74/40.40/43.05 | 50.66/45.96/47.83 | 38.74/35.73/36.80 |
| 手法② | 41.96/38.10/38.73 | 45.33/38.98/40.26 | 51.78/43.23/46.45 | 46.28/43.48/44.10 | 42.03/40.96/41.49 |
| 手法③ | 44.14/39.62/42.14 | 48.73/42.68/43.95 | 49.24/42.72/44.66 | 47.34/43.61/44.24 | 42.03/40.43/41.49 |

この結果を見る限り、手法③が他よりは良さそうに見える。しかしながら、LAS・MLAS・BLEX はいずれも 100 点満点の指標であり、その意味で表1は、全体に低調と言わざるを得ない。一体、どこに問題があるのか。

3 文切りの改善

手法①②③の処理結果(図4～6・24～26・44～46・64～66・84～86)を見る限り、筆者としては、どうも文切りがうまくいっていない印象を受ける。たとえば、2019年の手法①(図4)を正解(図22)と見比べてみると、正しく文が切り出せているのは「非敢当是也」「亦為報也」「県君有焉」の3文だけであり、他はズレてしまっている。文切りに失敗して、泣き別れになってしまった単語間のリンクは、必ず不正解になってしまうことから、これがLAS・MLAS・BLEX を引き下げてしまっているのだと考えられる。

この仮説を検証すべく、文切りを改善することを考えた。具体的には、センター試験の問題本文(図3・23・43・63・83)の句読点をもとに、文切り情報を事前に準備し、UDPipe の文切りに代えることとした。実際に比較したのは、以下の3つの手法である。

手法① 文切りを事前に与え、単語切り+品詞付与+係り受け解析に UDPipe を用いる

手法② 文切りを事前に与え、単語切りに MeCab を、品詞付与+係り受け解析に UDPipe を用いる

手法③ 文切りを事前に与え、単語切り+品詞付与に MeCab を、係り受け解析に UDPipe を用いる

処理結果を図7～9・27～29・47～49・67～69・87～89 に、比較結果を表2に示す。手法②と手法③が競っているものの、表1に比べると、LAS・MLAS・BLEX ともに改善され

ている。UDPipeの文切りオプション--tokenizer=joint_with_parsingは、非常に先進的な機能ではあるものの、現時点では、まだ試験的実装の域を出ていないことだろう。

表2: 手法①②③のLAS/MLAS/BLEX

| | 2019年 | 2018年 | 2017年 | 2016年 | 2015年 |
|-----|-------------------|-------------------|-------------------|-------------------|-------------------|
| 手法① | 71.58/67.72/70.89 | 51.10/45.40/47.85 | 60.05/50.49/54.37 | 59.47/53.87/56.97 | 58.11/55.70/56.23 |
| 手法② | 70.84/67.51/70.03 | 61.76/56.96/56.33 | 62.28/52.73/56.59 | 63.83/58.31/61.44 | 58.45/56.08/56.61 |
| 手法③ | 68.12/63.95/67.08 | 67.42/61.39/62.66 | 57.22/49.20/52.40 | 57.98/54.21/55.45 | 57.00/55.03/55.03 |

4 返り点をヒントにした品詞改善

センター試験の問題本文(図3・23・43・63・83)には、返り点が打たれている。これを元にして、手法①②③の処理結果(図7～9・27～29・47～49・67～69・87～89)を改善できないだろうか。

たとえば、2018年(図23)の5行目には、「所」←「沢」←「民」という返り点(上中下点)が打たれている。すなわち「沢」は、「民」からの返り先であると同時に、「所」への返り元である。古典中国語UDと返り点の関係(図2)を考慮すると、この場合、「沢」の品詞はVERB・AUX・PARTのいずれかにしか成り得ず、図27～29の「沢」の品詞NOUNは誤りである。NOUNは、返り元には現れ得るが、返り先には現れ得ないので、返り先かつ返り元である「沢」はNOUNでは有り得ない。

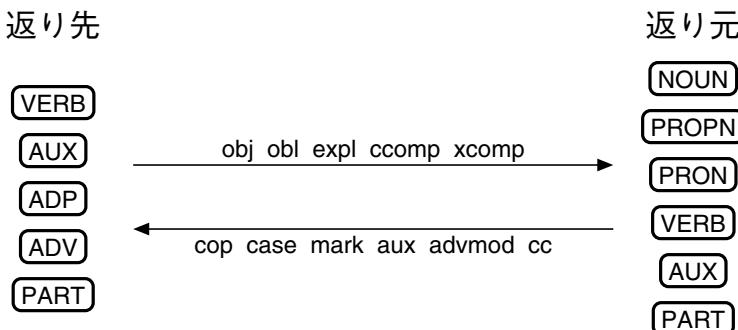


図2: 古典中国語UDと返り点の関係 [18]

このようなやり方で、問題本文中の全ての返り点をヒントに、手法①②③を、それぞれ改善するような手法を試した。

手法① 文切りを事前に与え、単語切り+品詞付与にUDPipeを用いた後、返り点による品詞改善をおこない、係り受け解析にUDPipeを用いる

手法② 文切りを事前に与え、単語切りにMeCabを、品詞付与にUDPipeを用いた後、返り点による品詞改善をおこない、係り受け解析にUDPipeを用いる

手法③ 文切りを事前に与え、単語切り+品詞付与に MeCab を用いた後、返り点による品詞改善をおこない、係り受け解析に UDPipe を用いる

処理結果を図 10~12・30~32・50~52・70~72・90~92 に、比較結果を表 3 に示す。2019 年・2018 年・2017 年は、表 2 より微妙な改善が見られるものの、2016 年と 2015 年は逆に悪化している。非常に残念だが、手法①②③は、かなりの作業量を必要とする割に、効果が薄い手法のようである。

表 3: 手法①②③の LAS/MLAS/BLEX

| | 2019 年 | 2018 年 | 2017 年 | 2016 年 | 2015 年 |
|-----|-------------------|-------------------|-------------------|-------------------|-------------------|
| 手法① | 72.68/69.21/72.38 | 55.49/48.93/52.60 | 64.12/55.19/59.74 | 48.42/41.49/42.72 | 55.69/51.46/54.11 |
| 手法② | 71.93/68.99/71.52 | 65.16/59.94/59.94 | 66.33/57.61/62.14 | 52.13/45.14/46.39 | 55.07/50.79/53.44 |
| 手法③ | 68.66/64.78/67.92 | 68.56/62.03/63.29 | 61.27/53.55/57.42 | 48.40/42.37/42.99 | 55.56/51.85/53.44 |

5 おわりに

大学入試センター試験の漢文を題材に、四書(孟子・論語・大学・中庸)を学んだ MeCab と UDPipe を用いて、古典中国語の単語切り・品詞付与・係り受け解析・文切りをおこなった。結論としては、単語切りに関しては MeCab の能力が高く、品詞付与では MeCab と UDPipe がトントンで、係り受け解析と文切りは UDPipe に任せることしかない。ただし、UDPipe の文切りは、古典中国語では良い結果が得られず、可能であれば、文切り情報を事前に準備した方がよい。一方、返り点をヒントにした品詞改善は、労おおくして功すべくなし、という感じである。

なお、本稿で学習に用いた古典中国語 UD コーパスや、MeCab・UDPipe などの学習済みモデルは

<http://kanji.zinbun.kyoto-u.ac.jp/~yasuoka/kyodokenkyu/2019-03-08/>

で公開している。古典中国語(漢文)に関する様々な用途に、どしどし使ってほしい。

付録 MeCab + StanfordNLP によるセンター漢文の解析

本稿の執筆中(2019 年 1 月 31 日)に、StanfordNLP 0.1.1 [24] がリリースされた。本稿の〆切まで時間が残されておらず、十分な解析は出来なかったものの、四書の学習、および MeCab との接続は実現できたので、以下の 9 つの手法を比較した。

手法Ⓐ 単語切り+品詞付与+係り受け解析+文切りの全てに StanfordNLP を用いる

手法Ⓑ 単語切りに MeCab を、品詞付与+係り受け解析+文切りに StanfordNLP を用いる

手法⑤ 単語切り + 品詞付与に MeCab を、係り受け解析 + 文切りに StanfordNLP を用いる

手法⑥ 文切りを事前に与え、単語切り + 品詞付与 + 係り受け解析に StanfordNLP を用いる

手法⑦ 文切りを事前に与え、単語切りに MeCab を、品詞付与 + 係り受け解析に Stanford NLP を用いる

手法⑧ 文切りを事前に与え、単語切り + 品詞付与に MeCab を、係り受け解析に Stanford NLP を用いる

手法⑨ 文切りを事前に与え、単語切り + 品詞付与に StanfordNLP を用いた後、返り点による品詞改善をおこない、係り受け解析に StanfordNLP を用いる

手法⑩ 文切りを事前に与え、単語切りに MeCab を、品詞付与に StanfordNLP を用いた後、返り点による品詞改善をおこない、係り受け解析に StanfordNLP を用いる

手法⑪ 文切りを事前に与え、単語切り + 品詞付与に MeCab を用いた後、返り点による品詞改善をおこない、係り受け解析に StanfordNLP を用いる

処理結果を図 13~21・33~41・53~61・73~81・93~101 に、比較結果を表 4 に示す。外部辞書を StanfordNLP に接続する方法を見つけられなかつたため、手法④⑤⑥⑦では、常用漢字などの品詞に間違いが多く、MLAS が非常に低くなっている。

表 4: 手法④⑤⑥⑦⑧⑨⑩⑪の LAS/MLAS/BLEX

| | 2019 年 | 2018 年 | 2017 年 | 2016 年 | 2015 年 |
|-----|-------------------|-------------------|-------------------|-------------------|-------------------|
| 手法④ | 31.40/28.12/28.12 | 39.34/33.54/34.16 | 43.04/32.90/36.13 | 32.28/27.27/28.48 | 36.01/32.88/31.27 |
| 手法⑤ | 33.79/29.02/30.28 | 42.49/34.50/36.42 | 43.04/32.90/36.13 | 34.57/28.75/30.58 | 35.27/32.09/30.48 |
| 手法⑥ | 40.33/36.19/38.10 | 43.06/35.03/37.58 | 46.08/38.83/40.78 | 42.55/37.42/39.26 | 43.48/41.69/41.69 |
| 手法⑦ | 46.70/38.34/40.26 | 51.38/43.48/45.34 | 48.61/37.54/40.78 | 39.68/33.33/33.33 | 45.26/41.71/41.18 |
| 手法⑧ | 47.41/39.24/40.51 | 60.62/50.32/55.41 | 48.61/37.54/40.78 | 42.02/34.89/36.14 | 45.41/40.96/40.96 |
| 手法⑨ | 62.67/55.70/60.13 | 66.29/59.94/58.68 | 54.68/47.74/50.32 | 51.60/45.34/47.83 | 52.66/49.74/50.79 |
| 手法⑩ | 50.00/42.81/42.81 | 51.38/43.48/45.34 | 49.11/38.19/40.78 | 41.27/34.46/35.08 | 45.26/41.71/41.18 |
| 手法⑪ | 50.68/43.67/43.04 | 61.76/51.59/56.69 | 49.11/38.19/40.78 | 43.62/36.02/37.89 | 45.41/40.96/40.96 |
| 手法⑫ | 62.67/55.70/60.13 | 66.86/60.57/59.31 | 55.70/49.03/51.61 | 51.06/45.34/47.83 | 52.66/49.74/50.79 |

また、ところどころにリンクの交差が散見される。われわれの古典中国語 UD では、リンクの交差は起こらない [15, 16] のだが、StanfordNLP は、リンクの交差をターゲットにしたアルゴリズム [25] を採用している。これが結果的にミスマッチとなっており、LAS・MLAS・BLEX を引き下げてしまっているように思われる。

参考文献

- [1] Joakim Nivre: Towards a Universal Grammar for Natural Language Processing, CICLing 2015: 16th International Conference on Intelligent Text Processing and Computational Linguistics (April 2015), pp.3-16.
- [2] Lucien Tesnière: *Éléments de Syntaxe Structurale*, Paris: C. Klincksieck (1959).
- [3] Igor A. Mel'čuk: Dependency Syntax: Theory and Practice, New York: State University of New York Press (1988).
- [4] 守岡知彦: MeCab を用いた古典中国語の形態素解析の試み, 情報処理学会研究報告, Vol.2008-CH-79 (2008 年 7 月), pp.17-22.
- [5] 守岡知彦: MeCab を用いた古典中国語形態素解析器の改良, 情報処理学会研究報告, Vol.2009-CH-84 (2009 年 10 月), No.3, pp.1-5.
- [6] Tomohiko Morioka: A Prototype of a Classical Chinese Morphological Analyzer based on MeCab, Proceedings of Osaka Symposium on Digital Humanities 2011 (September 2011), p.36.
- [7] 守岡知彦: 古典中国語形態素コーパス編集システムの開発, 東洋学へのコンピュータ利用, 第 23 回研究セミナー (2012 年 3 月), pp.75-83.
- [8] 山崎直樹, 守岡知彦, 安岡孝一: 古典中国語形態素解析のための品詞体系再構築, 人文科学とコンピュータシンポジウム「じんもんこん 2012」論文集 (2012 年 11 月), pp.39-46.
- [9] Tomohiko Morioka, Christian Wittern, Koichi Yasuoka, Naoki Yamazaki: A Study of Linguistic Analysis for Classical Chinese Texts, Proceedings 2013 International Conference on Culture and Computing (September 2013), pp.143-144.
- [10] Koichi Yasuoka, Naoki Yamazaki, Christian Wittern, Yoshihiro Nikaido, Tomohiko Morioka: A Morphological Analysis of Classical Chinese Texts, Proceedings of Digital Humanities 2014 (July 2014), pp.410-412.
- [11] 安岡孝一, 守岡知彦, Christian Wittern, 山崎直樹, 二階堂善弘, 鈴木慎吾: 古典中国語形態素解析による地名の自動抽出, 人文科学とコンピュータシンポジウム「じんもんこん 2014」論文集 (2014 年 12 月), pp.63-68.
- [12] 安岡孝一, Christian Wittern, 守岡知彦, 池田巧, 山崎直樹, 二階堂善弘, 鈴木慎吾, 師茂樹: 古典中国語(漢文)の形態素解析, 東洋学へのコンピュータ利用, 第 27 回研究セミナー (2016 年 3 月), pp.3-14.
- [13] 安岡孝一, ウィッテルン クリストイアン, 守岡知彦, 池田巧, 山崎直樹, 二階堂善弘, 鈴木慎吾, 師茂樹: 古典中国語(漢文)の形態素解析とその応用, 情報処理学会論文誌, Vol.59, No.2 (2018 年 2 月), pp.323-331.

- [14] 安岡孝一, ウィッテルンクリスティアン, 守岡知彦, 池田巧, 山崎直樹, 二階堂善弘, 鈴木慎吾, 師茂樹: 古典中国語 Universal Dependencies への挑戦, 情報処理学会研究報告, Vol.2018-CH-116 (2018年1月), No.20, pp.1-8.
- [15] 守岡知彦: 古典中国語 UD コーパスの IPFS を用いた表現の試み, 情報処理学会研究報告, Vol.2018-CH-118 (2018年8月), No.6, pp.1-7.
- [16] 安岡孝一: 古典中国語(漢文)の依存文法解析と直接構成素解析, 漢字文献情報処理研究, 第18号(2018年10月), pp.56-62.
- [17] 安岡孝一: Universal Dependencies にもとづく古典中国語(漢文)の依存文法解析, センター研究年報2018(2018年10月).
- [18] 安岡孝一: 漢文の依存文法解析と返り点の関係について, 日本漢字学会第1回研究大会予稿集(2018年12月), pp.33-48.
- [19] 安岡孝一: 古典中国語 Universal Dependencies で読む『孟子』, センター研究年報2018別冊(2019年3月).
- [20] Taku Kudo, Kaoru Yamamoto, Yuji Matsumoto: Applying Conditional Random Fields to Japanese Morphological Analysis, Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (July 2004), pp.230-237.
- [21] Milan Straka and Jana Straková: Tokenizing, POS Tagging, Lemmatizing and Parsing UD 2.0 with UDPipe, Proceedings of the CoNLL 2017 Shared Task (August 2017), pp.88-99.
- [22] ウィッテルン・クリスティアン: 漢籍リポジトリ, センター研究年報2015(2016年3月).
- [23] Daniel Zeman, Jan Hajič, Martin Popel, Martin Potthast, Milan Straka, Filip Ginter, Joakim Nivre, and Slav Petrov: CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies, Proceedings of the CoNLL 2018 Shared Task (October 2018), pp.1-21.
- [24] Peng Qi, Timothy Dozat, Yuhao Zhang, Christopher D. Manning: Universal Dependency Parsing from Scratch, Proceedings of the CoNLL 2018 Shared Task (October 2018), pp.160-170.
- [25] Timothy Dozat, Peng Qi, Christopher D. Manning: Stanford's Graph-based Neural Dependency Parser at the CoNLL 2017 Shared Task, Proceedings of the CoNLL 2017 Shared Task (August 2017), pp.20-30.

鳴呼哀哉。有_二兄_一子_一曰_フ甫_ト制_シ服_ヲ於_ニ斯_ニ紀_ニ德_ヲ於_ニ斯_ニ刻_ニ石_ニ於_ニ斯_ニ
(注1) (注2)

或_一曰_{ハク}「豈_ニ孝童之猶子_{ナル}与_カ奚_{なん}孝義之勤若此。」甫泣_{キテ}而對_(ア)曰_{ハク}「非_ニザル」
(注4)

敢_{ヘテ}當_{タルニ}是_{レニ}也_タ亦_ス為_ル報_{ユルヲ}也_。甫昔臥_{フシ}病_ニ於_ガ我_諸姑_ニ姑_ニ之_子又_ム病_ム問_{ハバ}
(注5)

女_一巫_ニ巫_一曰_{ハク}「廻_シ檻_ノ之_東南_隅者_吉。」姑遂_{カヘ}易_シ予_ヲ之_地以_テ安_{シズ}我_ヲ我_ヲ
(注6)

將_レ出_レ涕_{ダサ}、存_シ而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注7)

用_レ是_{レヲ}而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注8)

将_{マサニ}出_レ涕_{ダサ}、存_シ而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注9)

用_レ是_{レヲ}而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注10)

用_レ是_{レヲ}而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注11)

用_レ是_{レヲ}而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注12)

用_レ是_{レヲ}而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注13)

用_レ是_{レヲ}而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注14)

用_レ是_{レヲ}而_{シテ}姑_ニ之_子卒_シ。後_ニ乃_レ知_ル之_ヲ於_ガ走_シ使_{ヨリ}甫嘗_{カツテ}有_リ說_{クコト}於_テ人_ニ客_ニ
(注15)

A 甫_ト制_シ服_ヲ於_ニ斯_ニ紀_ニ德_ヲ於_ニ斯_ニ刻_ニ石_ニ於_ニ斯_ニ
 B 非_ニザル

C 小括

D 安_{シズ}我_ヲ我_ヲ

E 小括

F 小括

G 小括

H 小括

I 小括

J 小括

K 小括

L 小括

M 小括

N 小括

O 小括

P 小括

Q 小括

R 小括

S 小括

T 小括

U 小括

V 小括

W 小括

X 小括

Y 小括

Z 小括

図 3: 大学入試センター試験『国語』(2019年1月19日)第4問本文

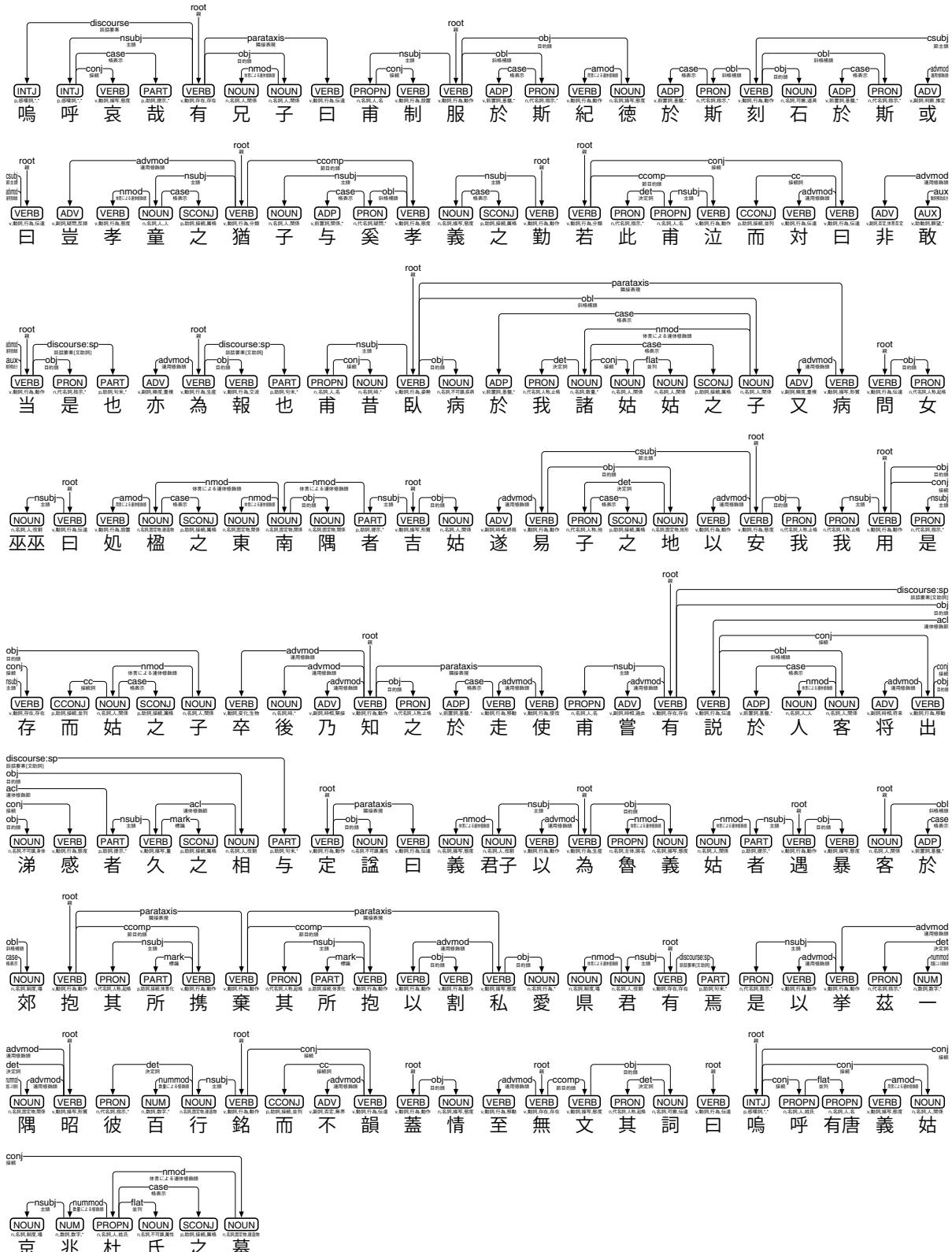


図 4: 手法①の処理結果 (2019 年)

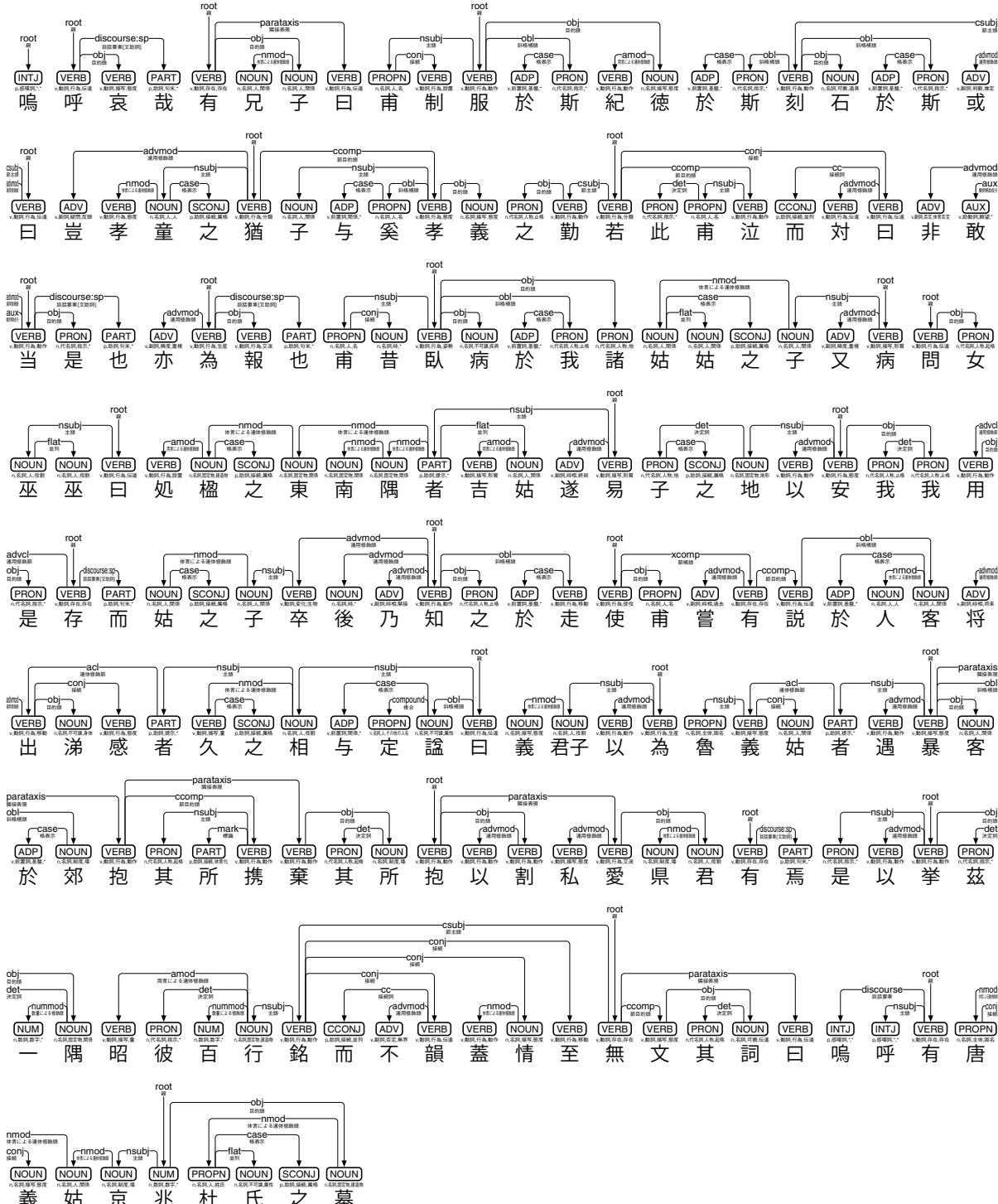


図 5: 手法②の処理結果(2019 年)

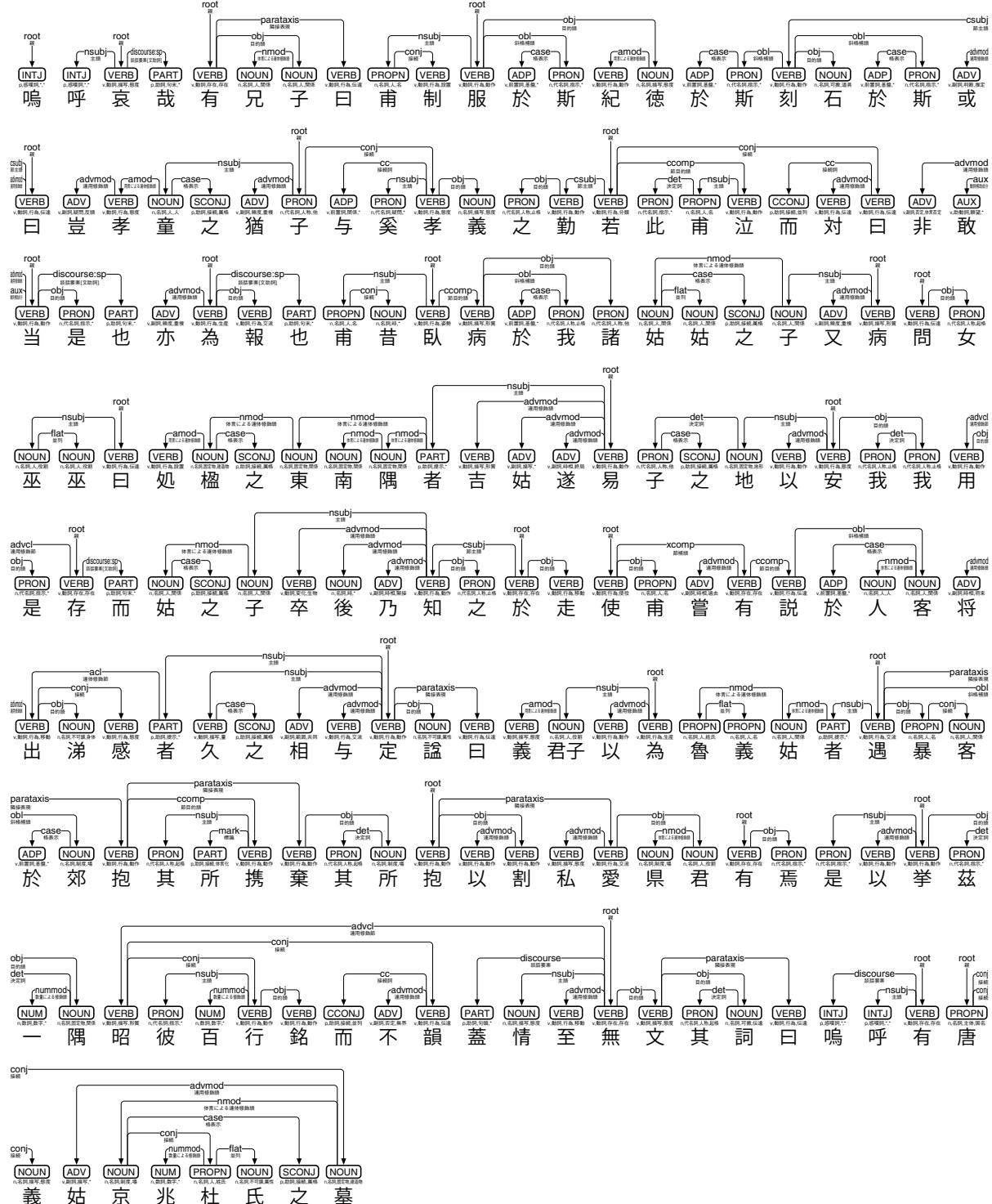


図 6: 手法③の処理結果(2019年)

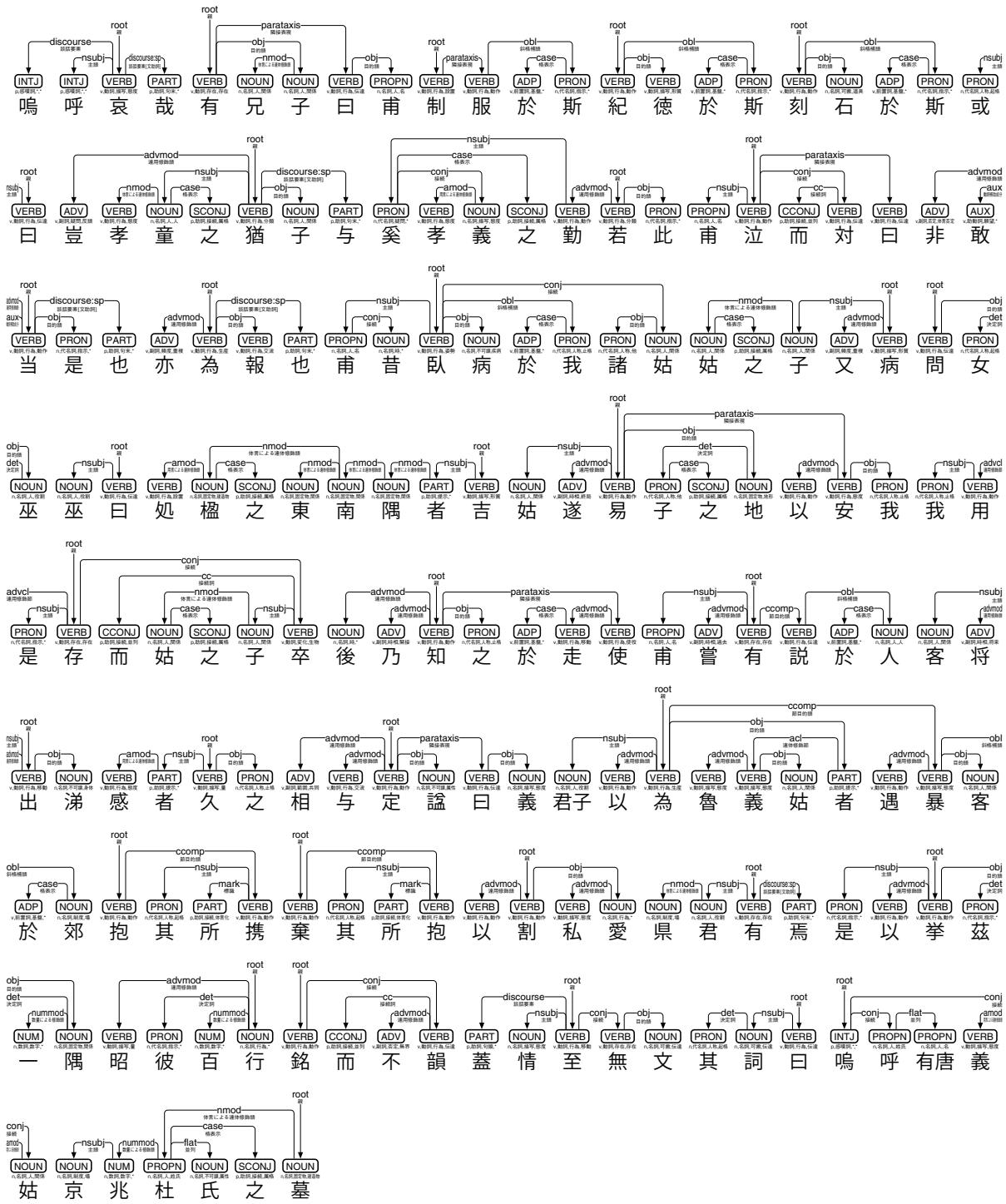


図 7: 手法①の処理結果(2019年)

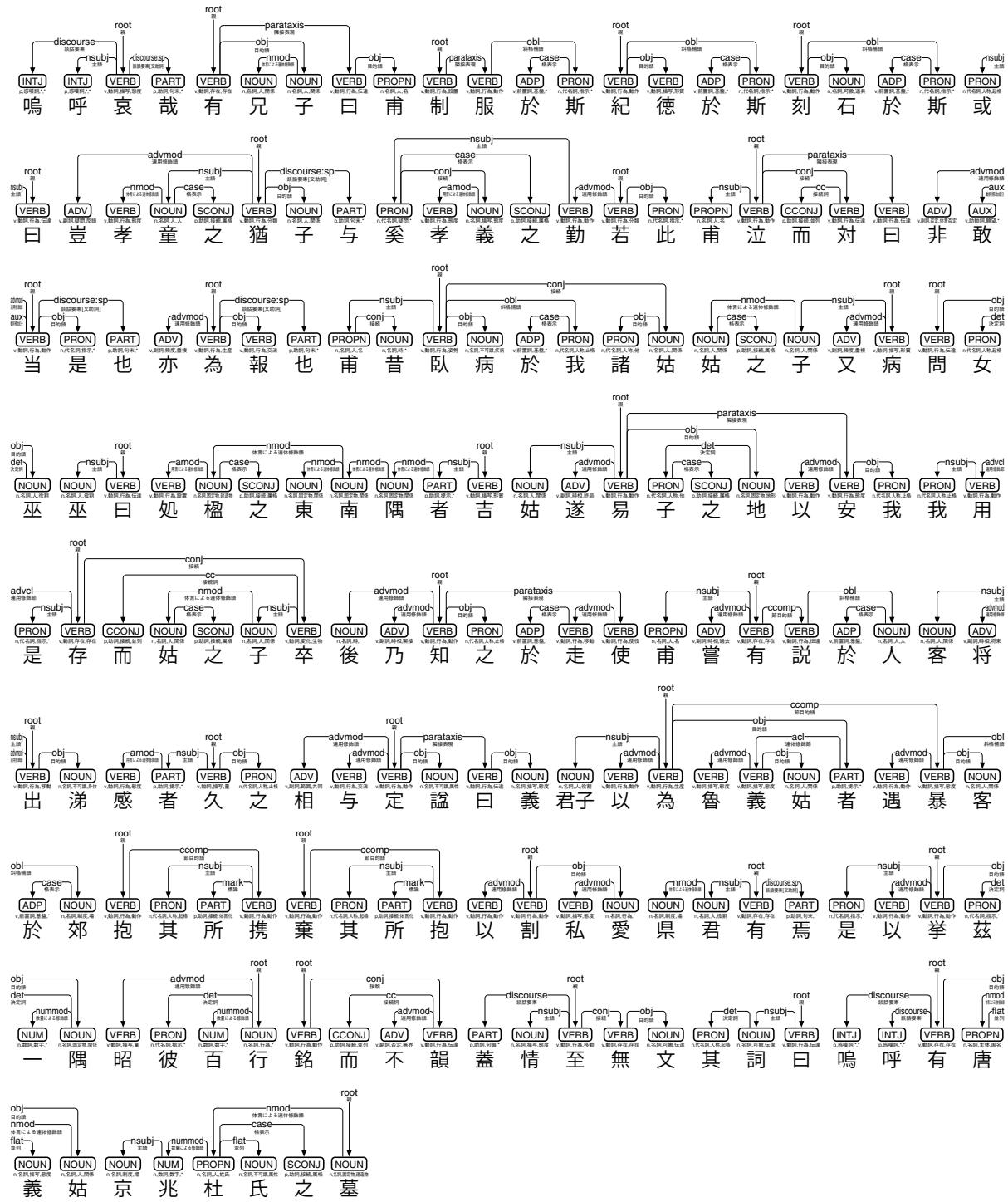


図 8: 手法②の処理結果(2019年)

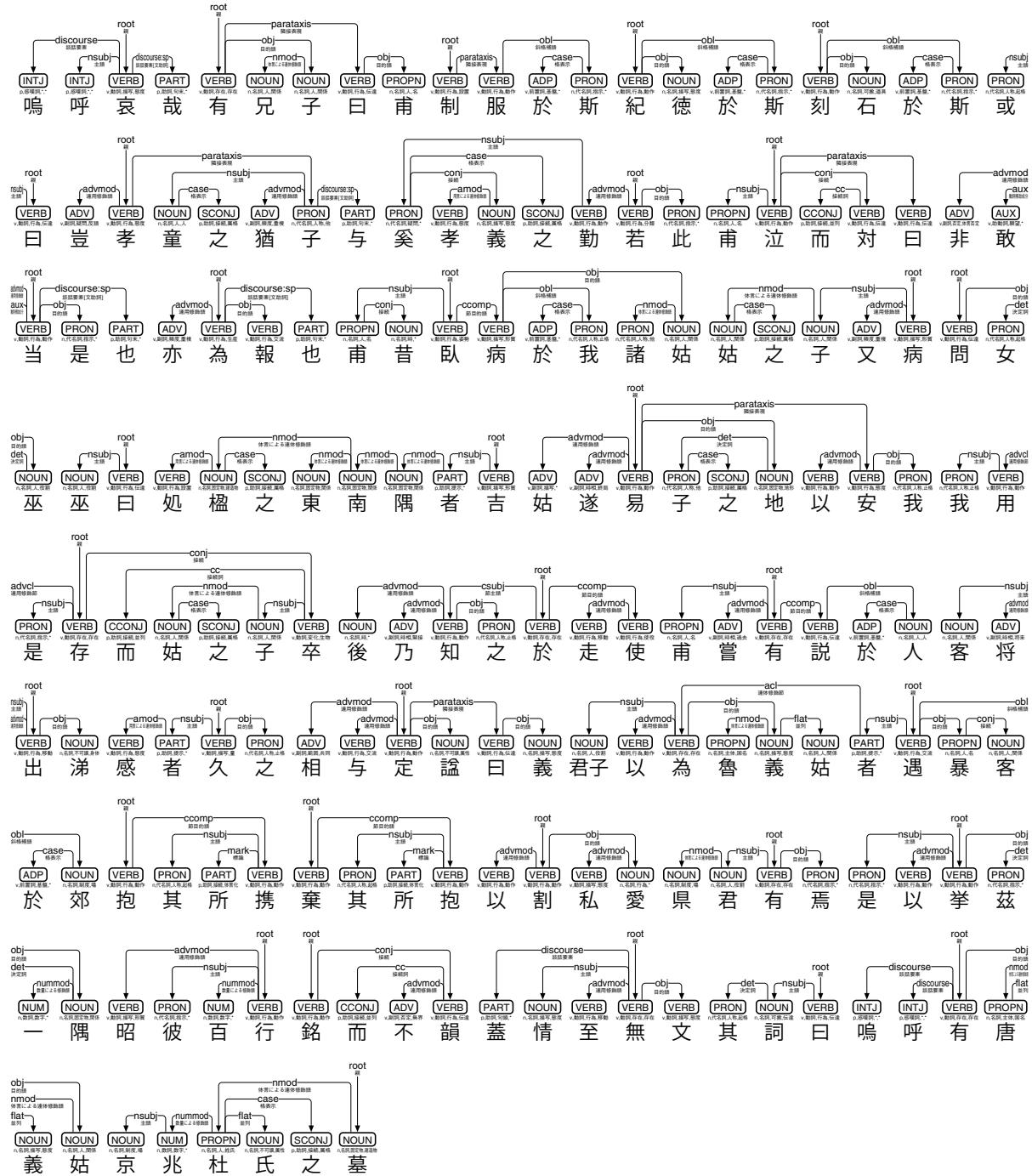


図9: 手法③の処理結果(2019年)

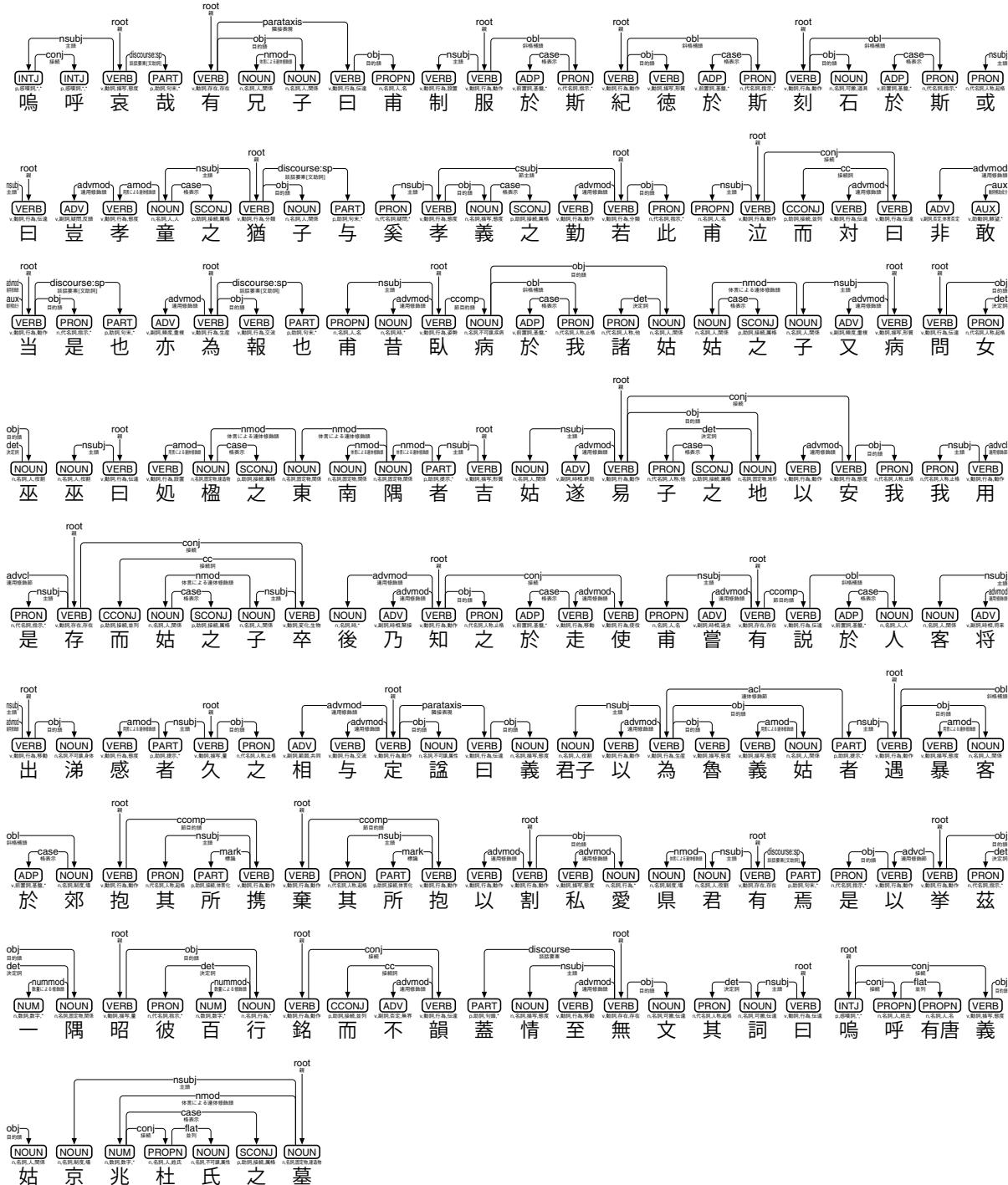


図 10: 手法①の処理結果(2019 年)

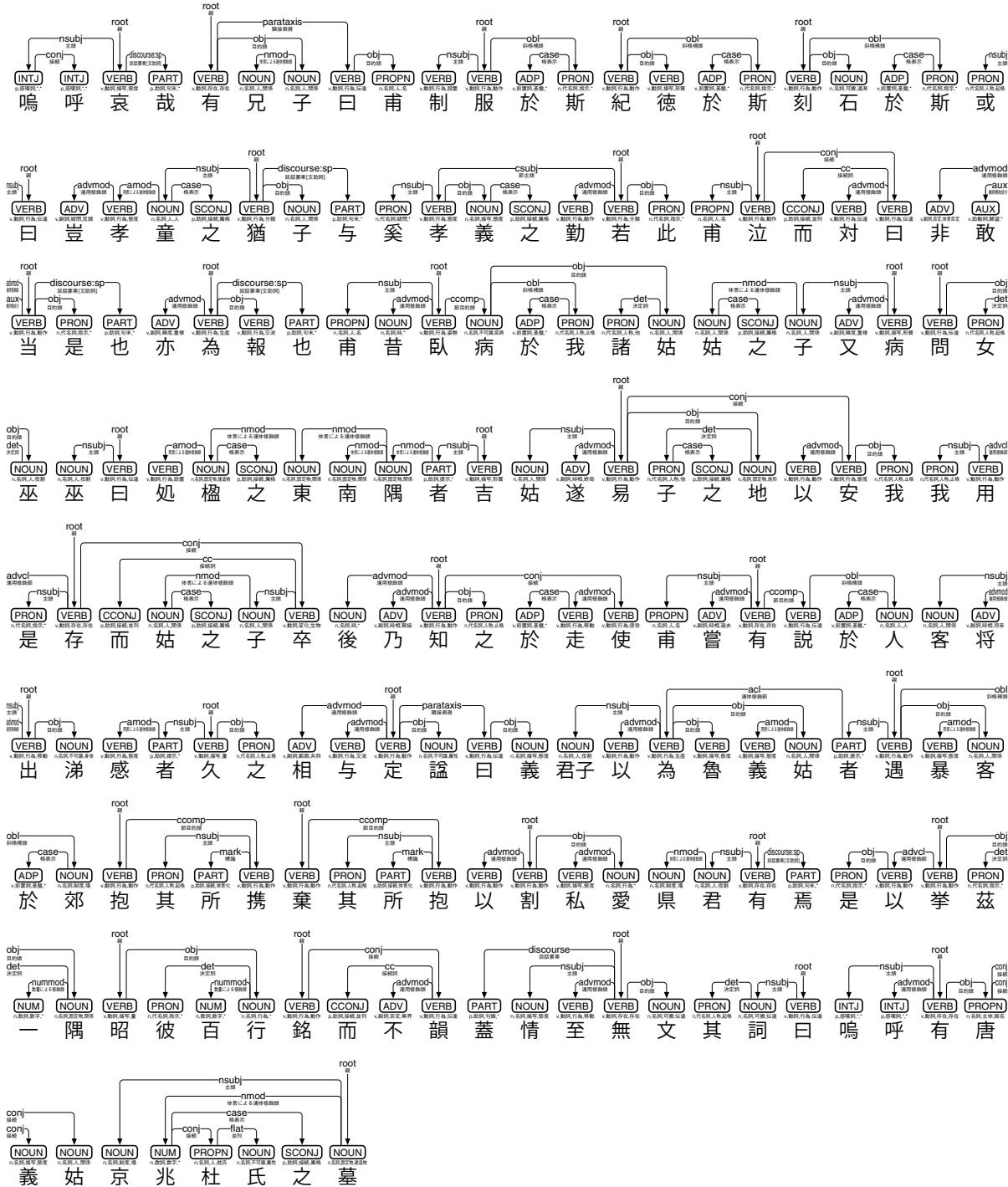


図 11: 手法②の処理結果(2019年)

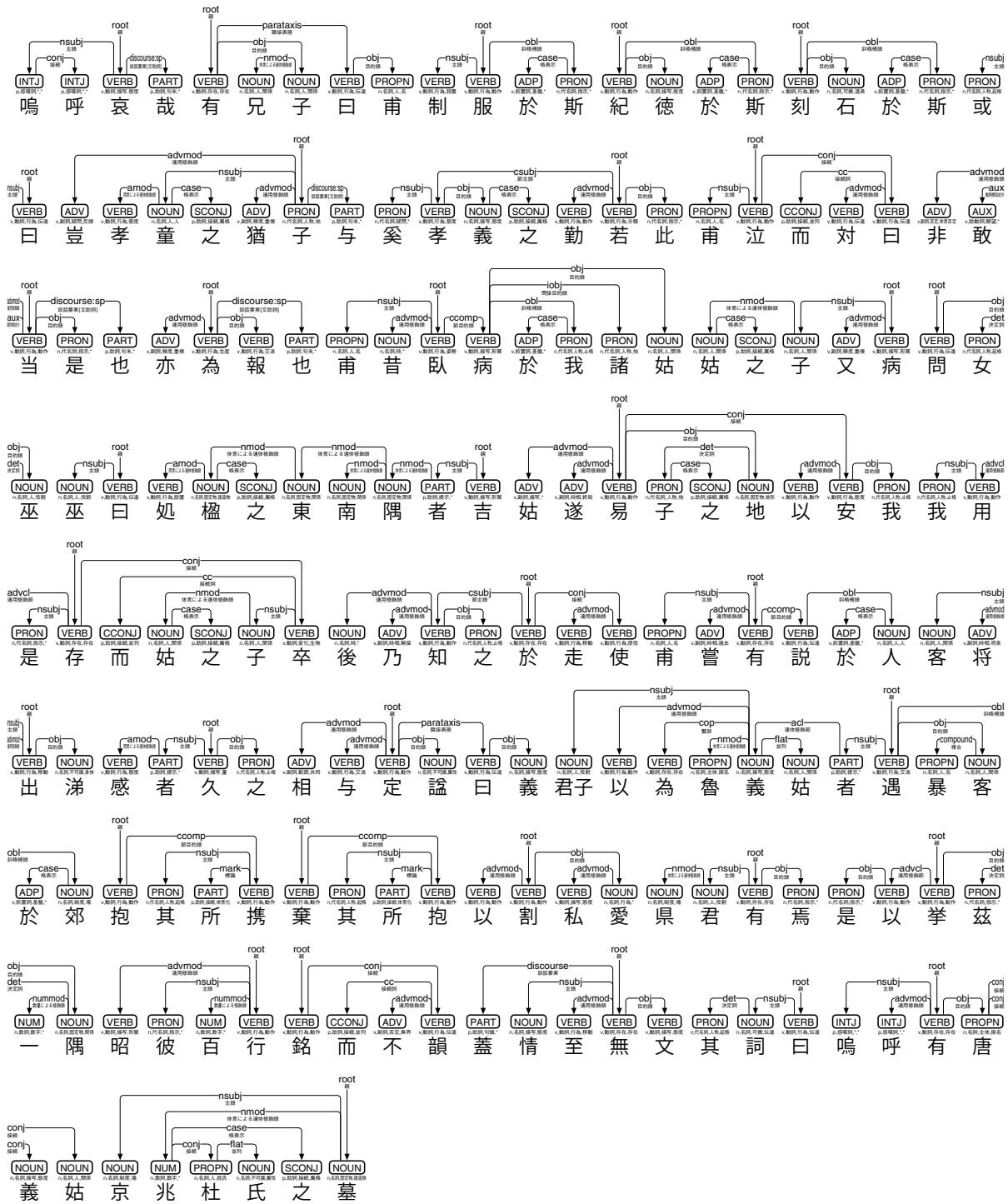


図 12: 手法③の処理結果(2019年)

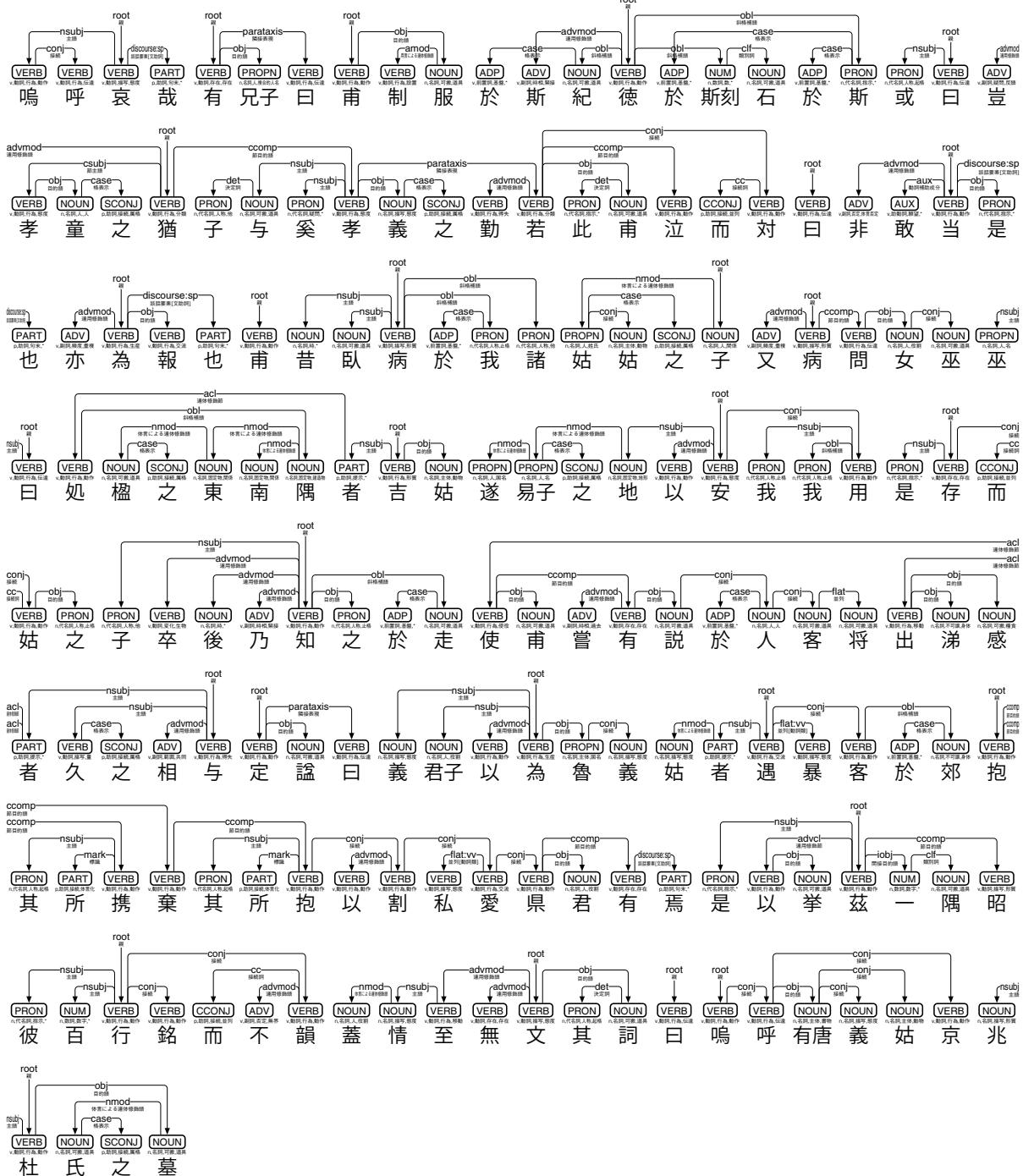


図 13: 手法Ⓐの処理結果(2019年)

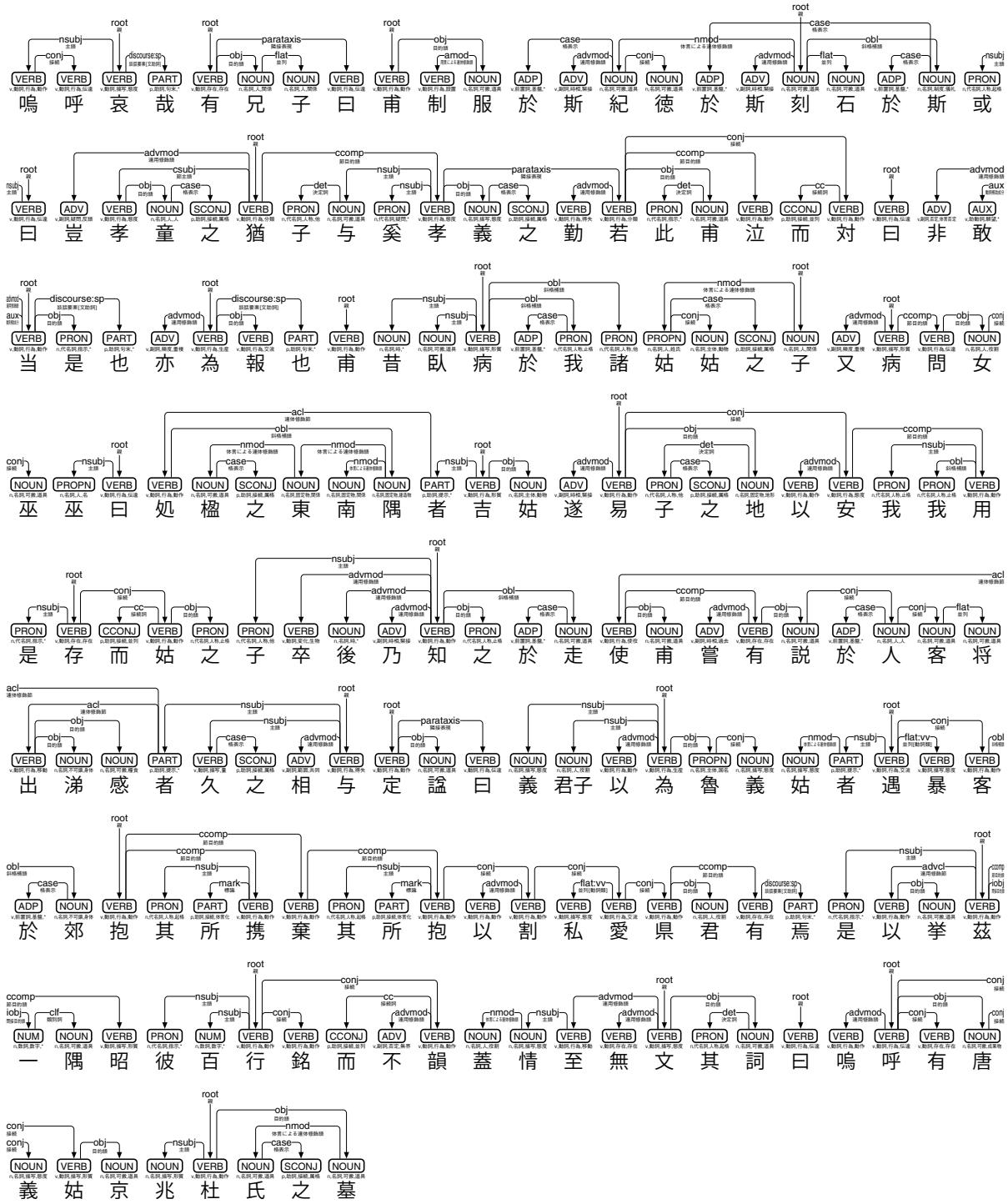


図 14: 手法Bの処理結果(2019年)

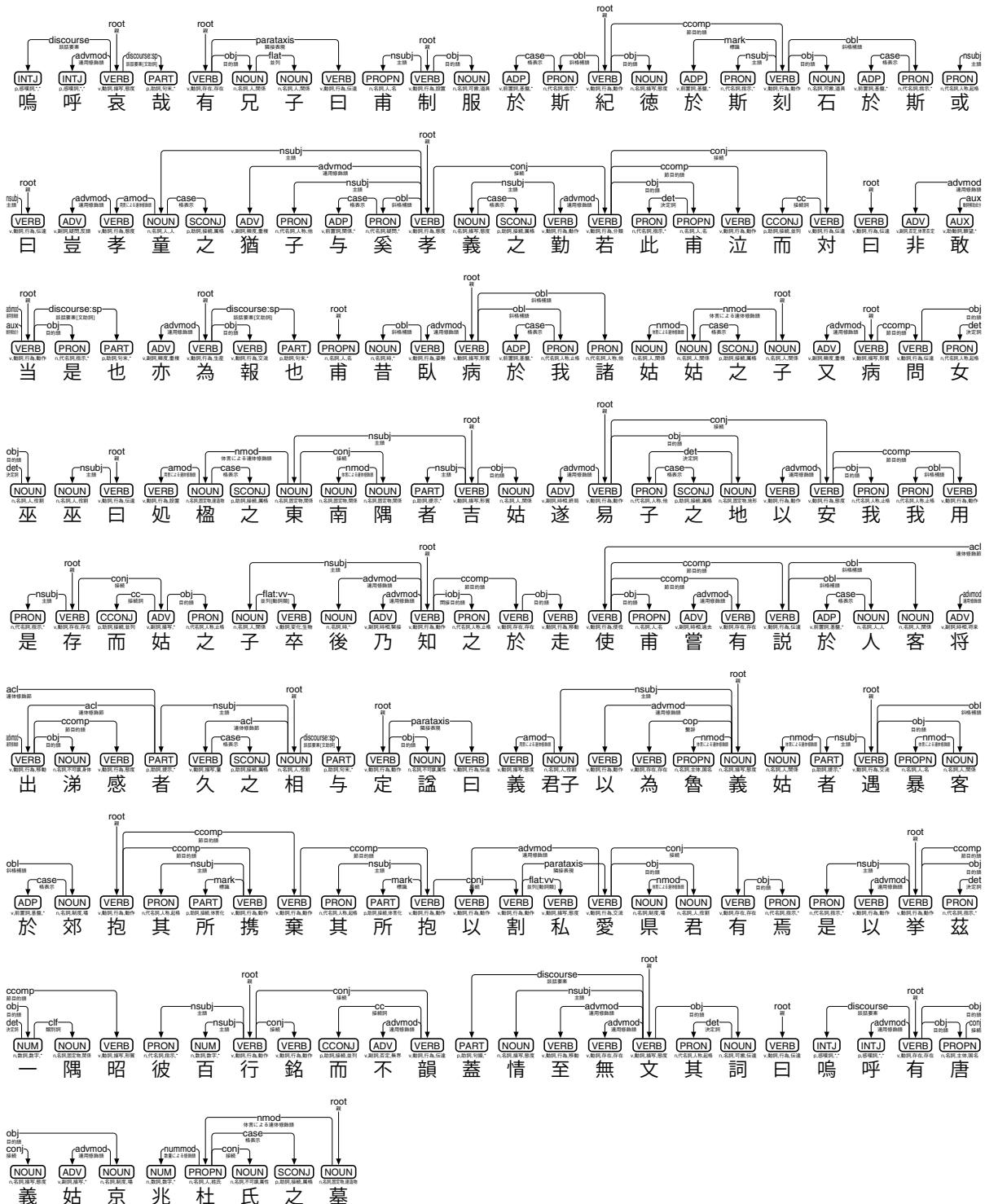


図 15: 手法④の処理結果(2019年)

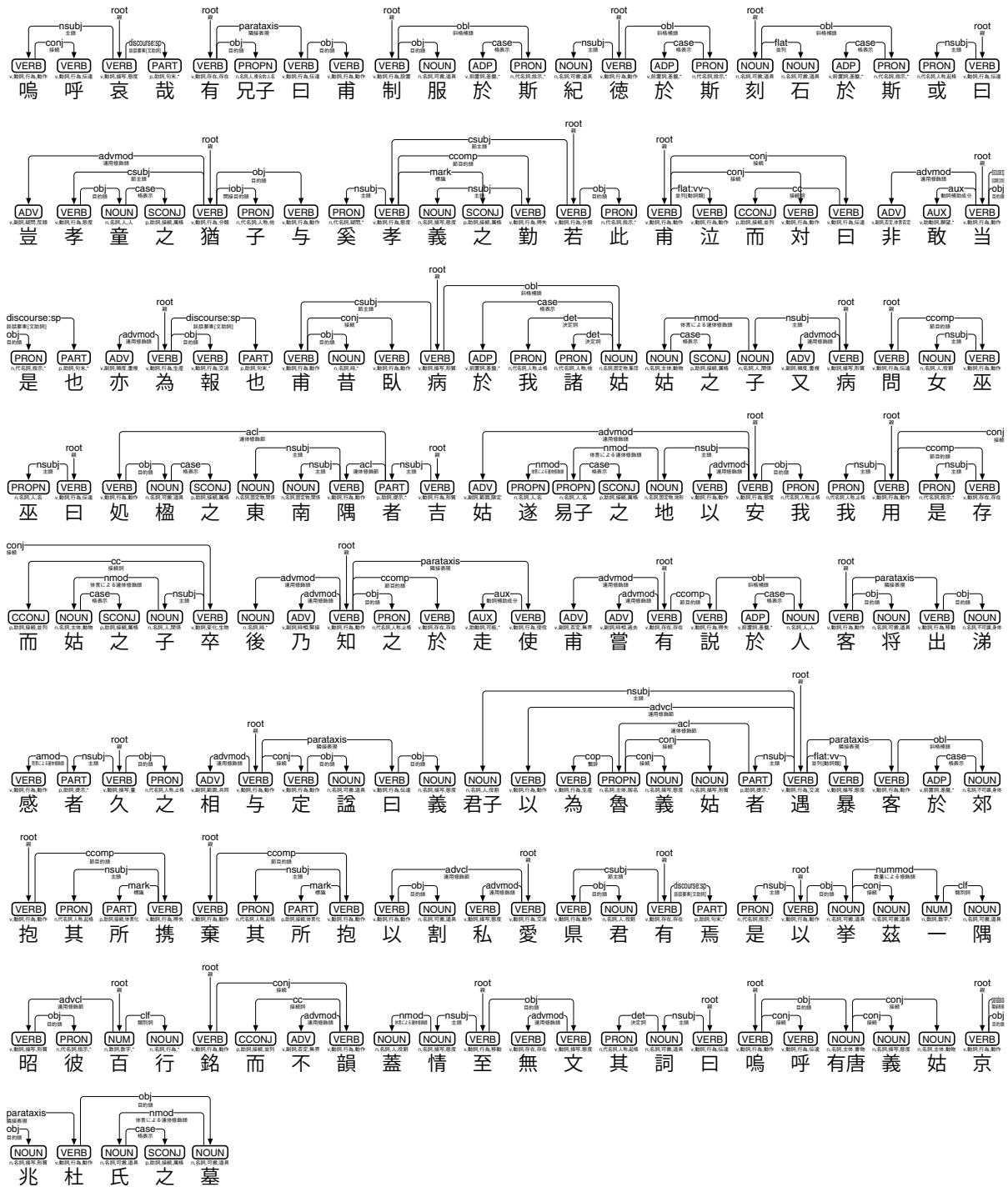


図 16: 手法Aの処理結果(2019年)

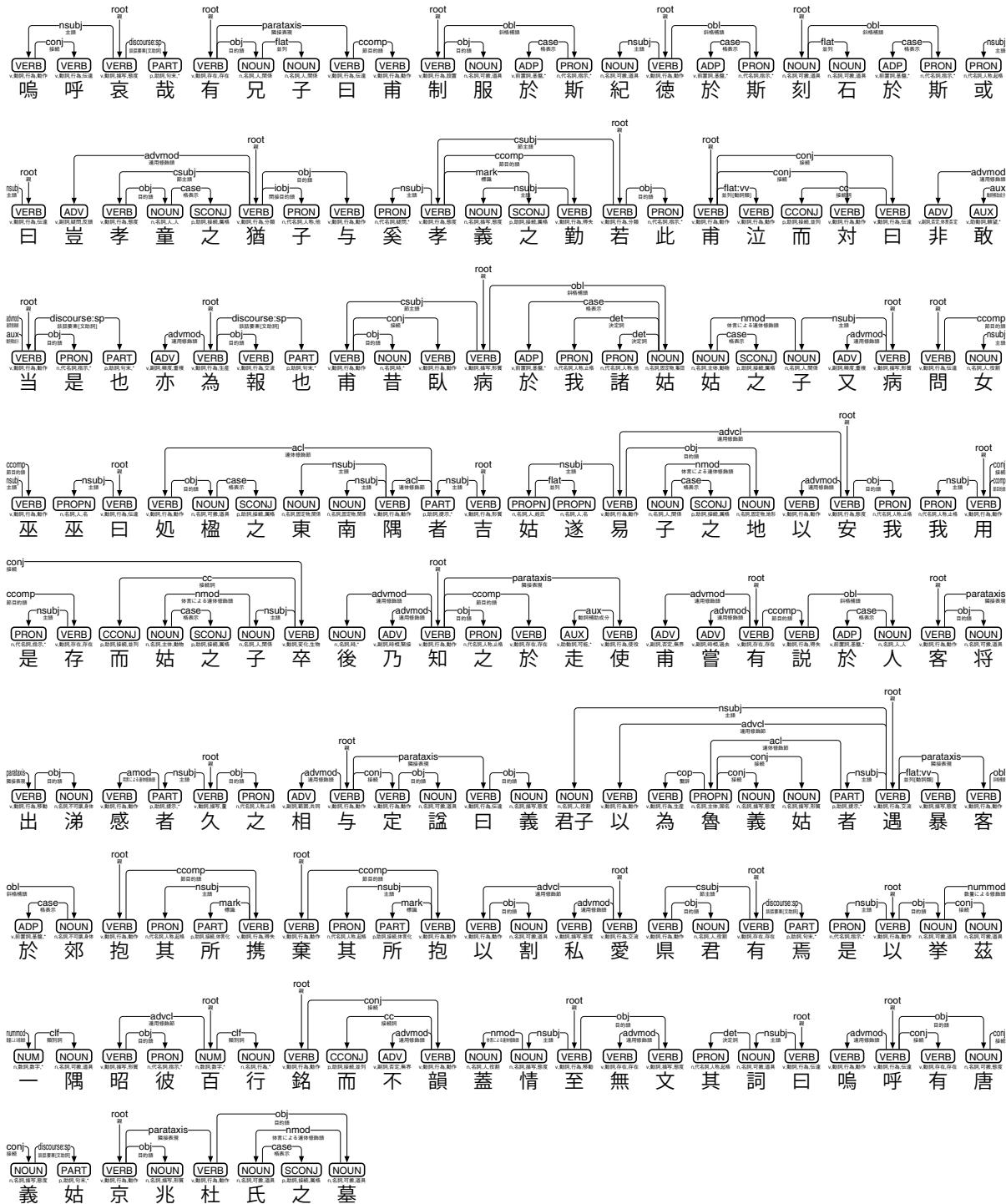


図 17: 手法Bの処理結果(2019年)

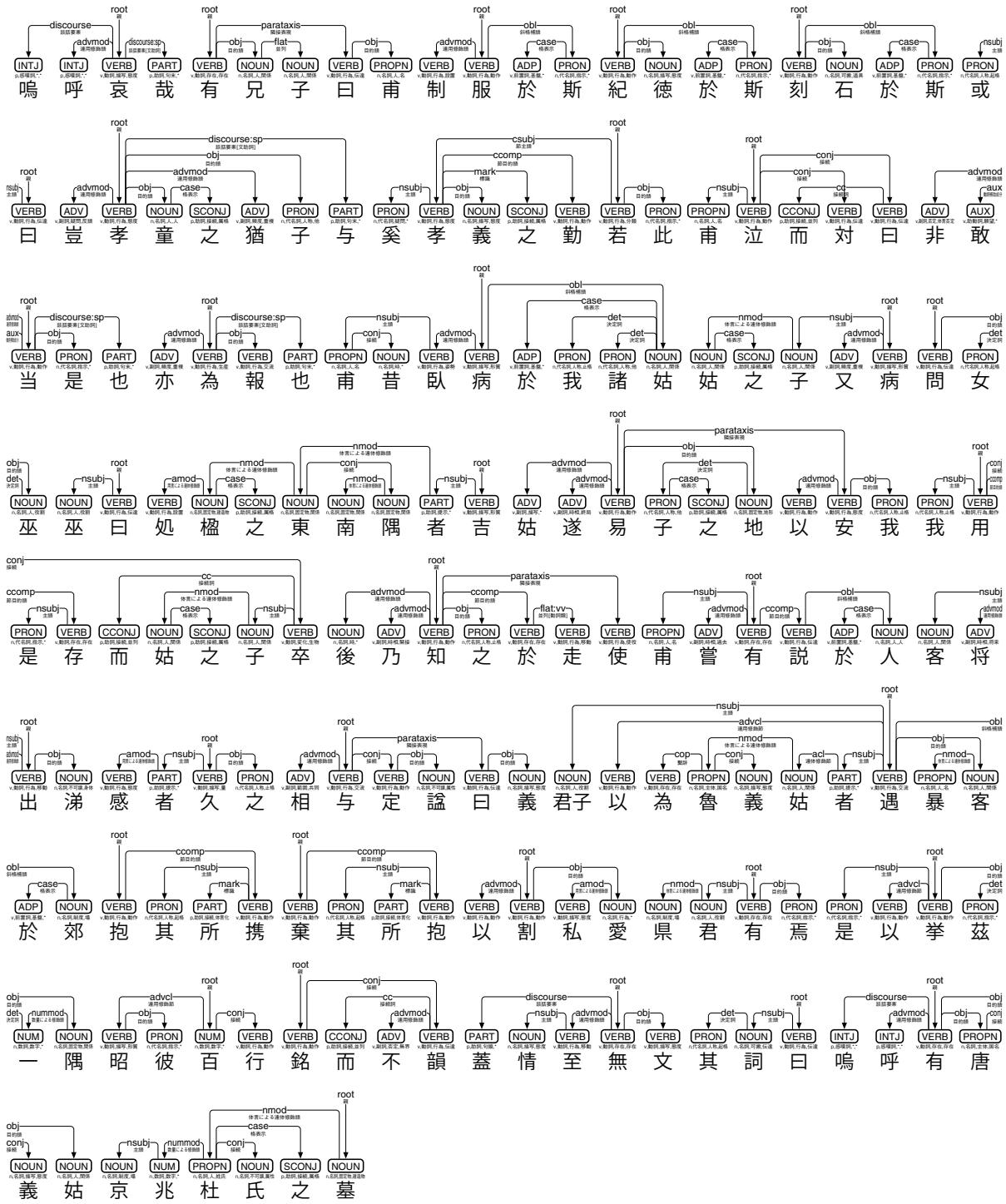


図 18: 手法Cの処理結果(2019年)

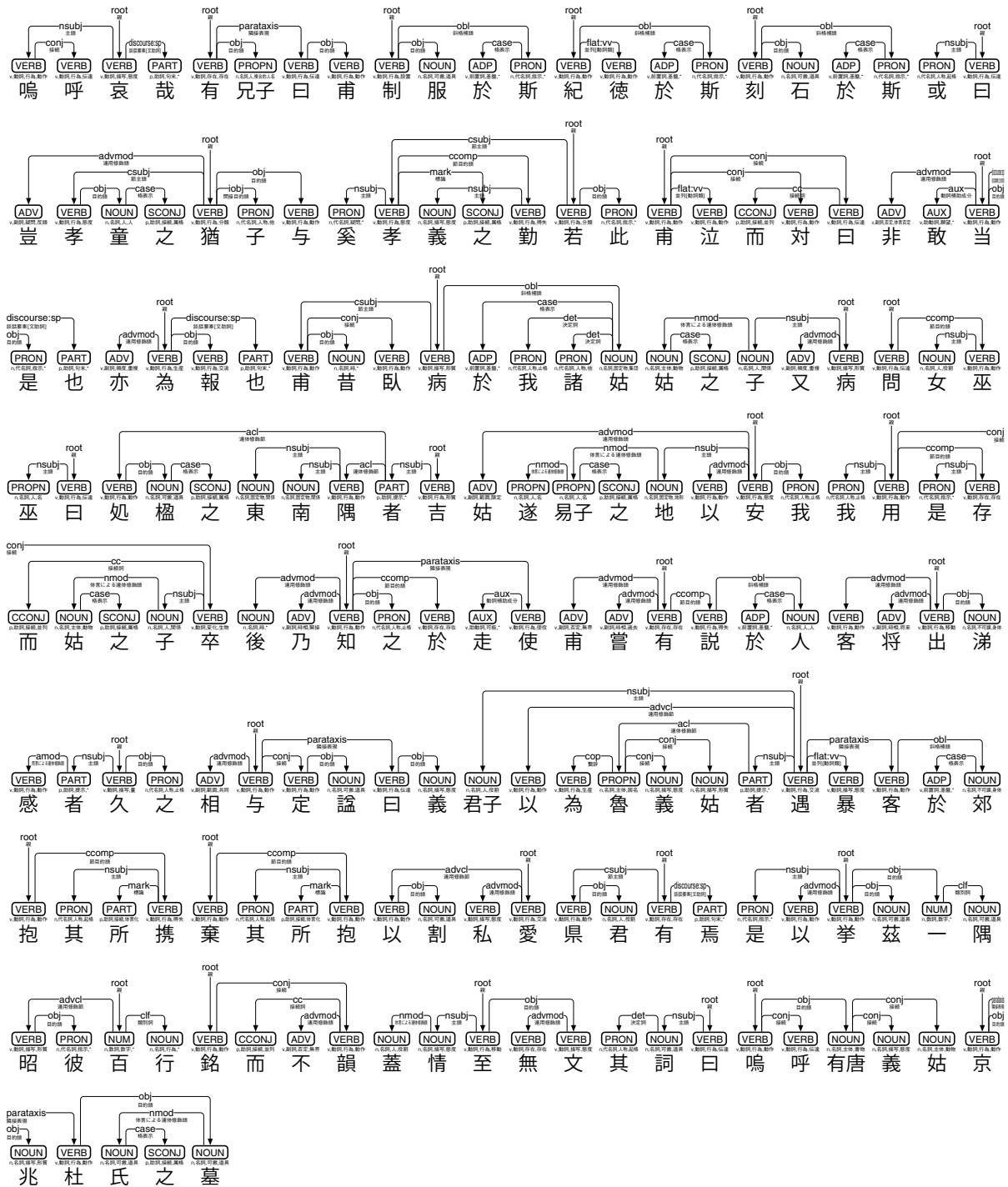


図 19: 手法Aの処理結果(2019年)

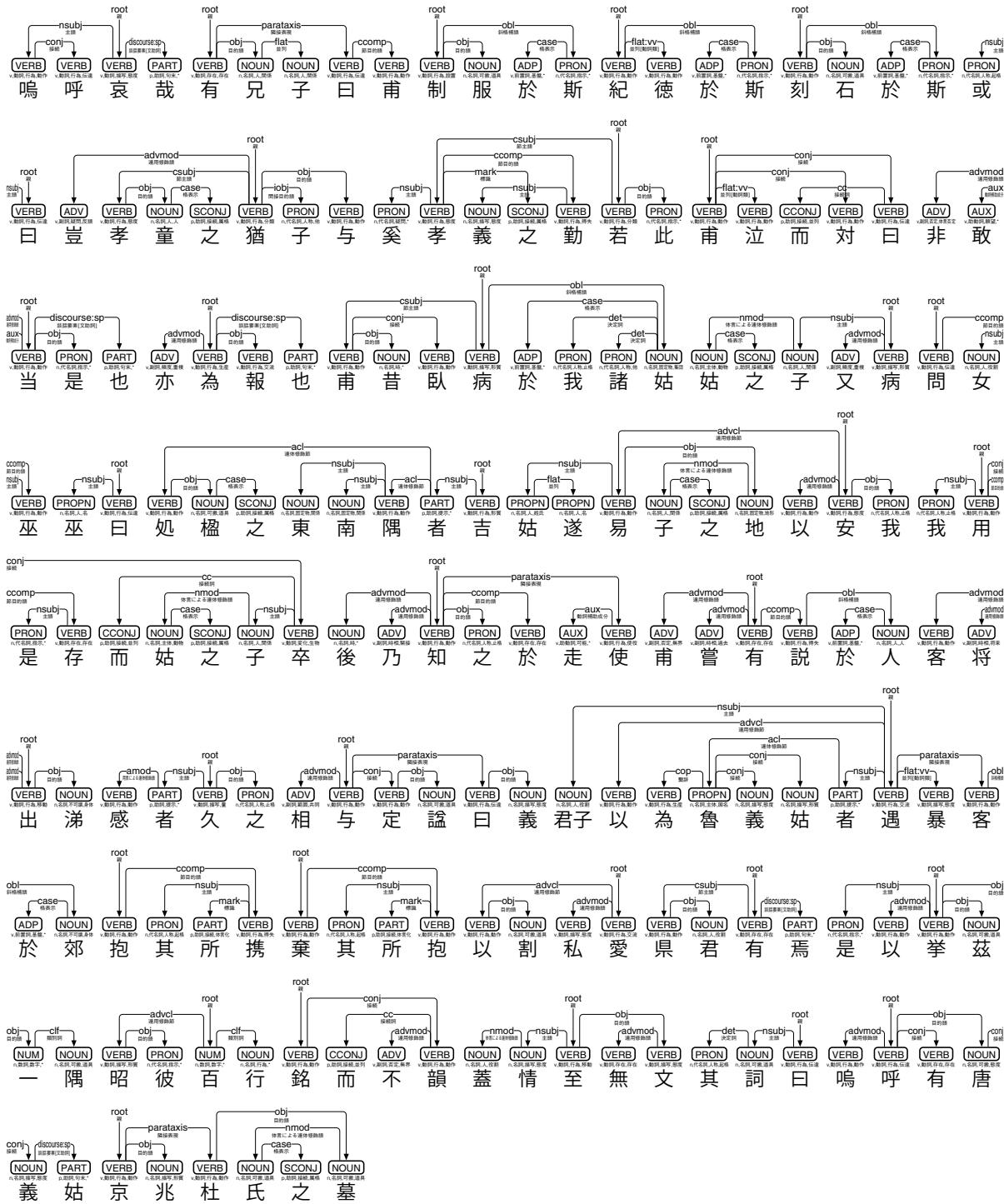


図 20: 手法Bの処理結果(2019年)

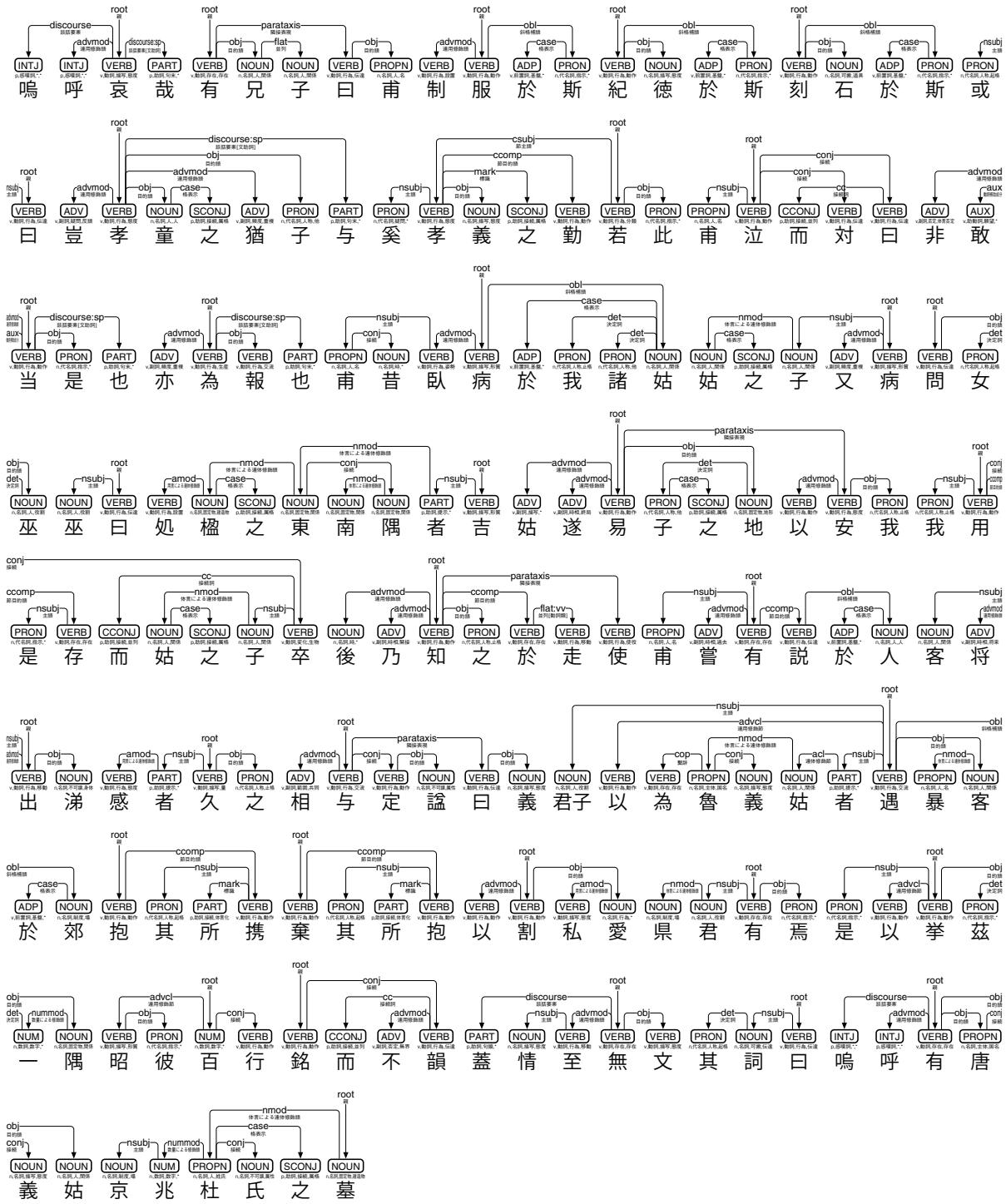


図 21: 手法Cの処理結果(2019年)

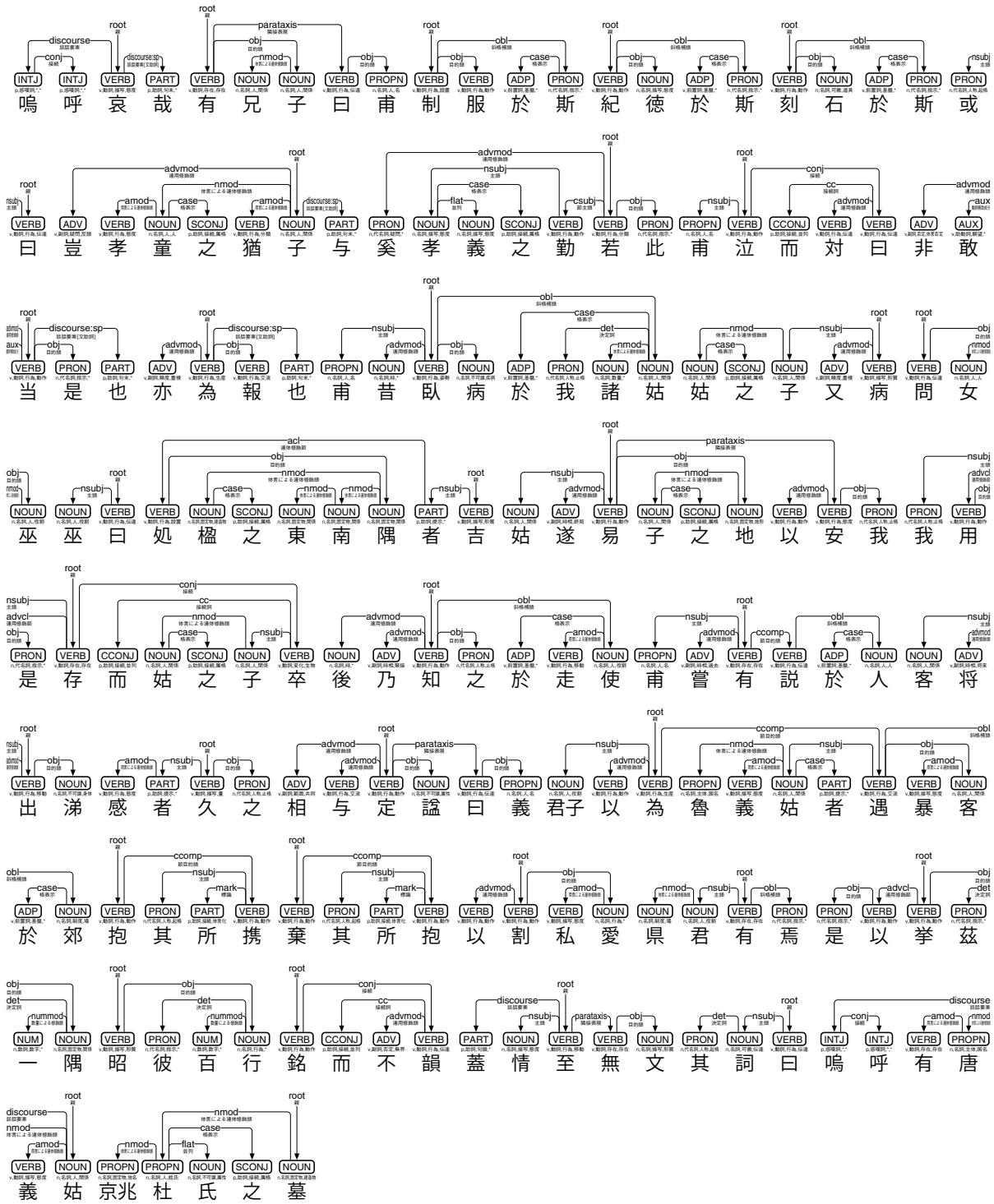


図 22: 手作業で作成した「正解」UD (2019 年)

| | | | | | | | |
|-----------------------------|---|--|--|---|---|---|---|
| 於深識遠慮殆不 _レ 能勝吾子也。 | D 恐 _ル 譽望之損 _一 也。準喜、起執其手 _一 曰、元 _{げん} 之雖 _モ 文 _一 章冠 _{タリト} 天 _一 下 _ニ 至 _{リテハ} | B 以 _テ 太平 _ヲ 責 _{もとメン} 焉。丈人之于 _ニ 明 _{おケルヤ} ^(注14) 主、能 _ク 若 _キ 魚 _一 之有 _{ルガ} 水 _{乎。} 嘉祐所 _ニ 以 | C 故言聽計徒而功名俱美。今丈人負 _ニ 天 _一 下 _ノ 重 _ニ 相 _{タレバ} 則 _チ 中 _外 | Y 相所以能建 _ニ 功業 _ヲ 澤 _ニ 生 _メ 民 _上 者、其君臣相 _ヒ 得皆如 _ニ 魚 _一 之有 _{ルガ} 水 _。 | X 不若未為相。為相則 _ハ 譽望損矣。準曰、何故 _{ゾト} 嘉祐曰、自古賢 | Z 封府一日、問 _ニ 嘉祐 _一 曰、外間議 _レ 準 _ヲ 云 _カ 何 _。 嘉祐曰、外人皆云 _ニ 丈 _フ 人 _ト | I 嘉祐禹偁子也。嘉祐平時若愚騃 _{カント} 獨 _リ 寇 _{コラ} 準 _{ジュン} 知 _レ 之。準知 _ニ 開 _{ミシカ} |
| | | | | | A 曰 _{ハク} 於 _ニ 吾 _子 意 _ヲ 何 _。 如 _{ハク} 嘉祐 _一 曰、愚 _ヲ 觀 _レ 之、丈人 _。 | | |
| | | | | | II <small>注5</small> | | |

図 23: 大学入試センター試験『国語』(2018年1月13日)第4問本文

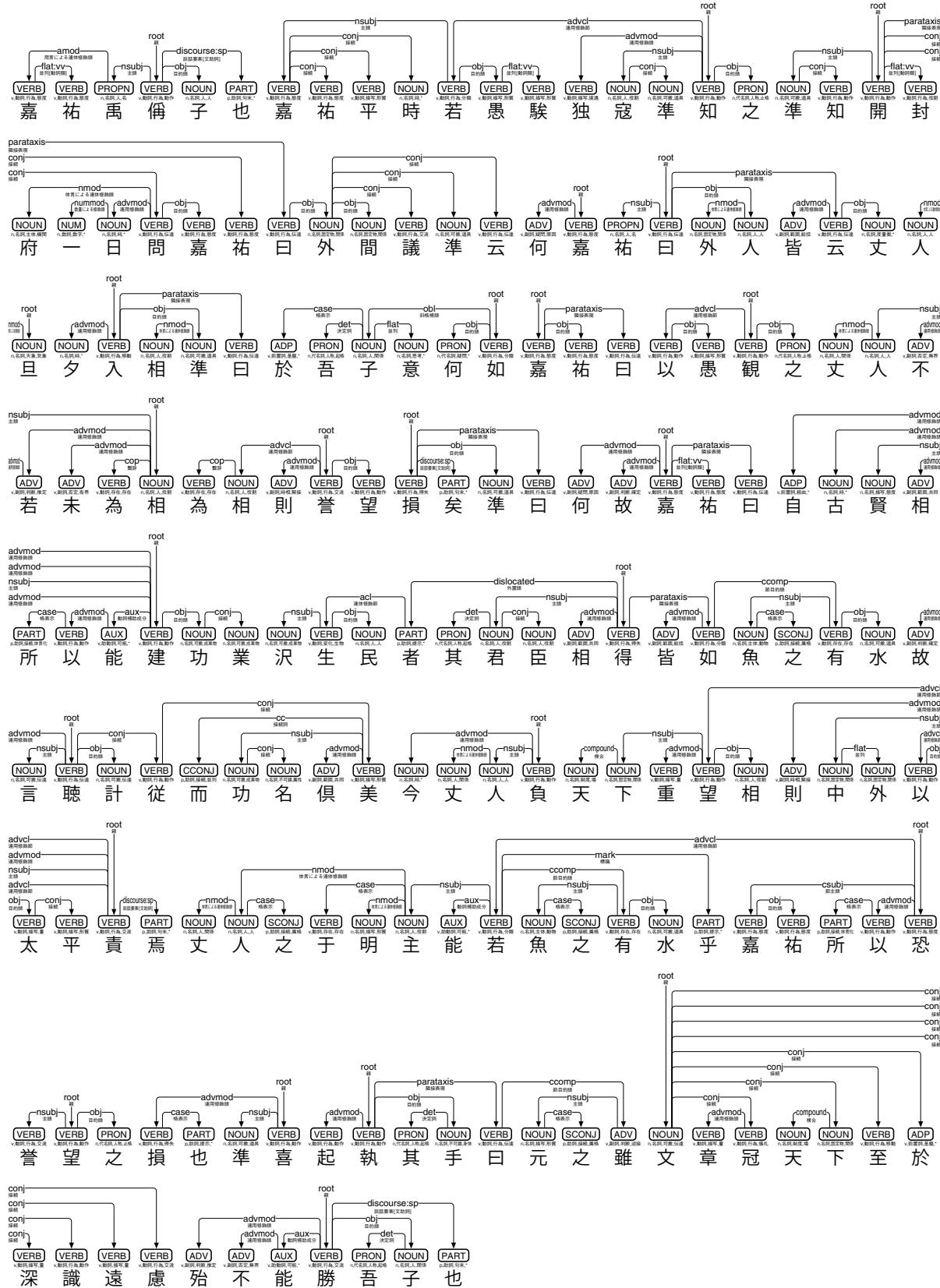


図 24: 手法①の処理結果(2018年)

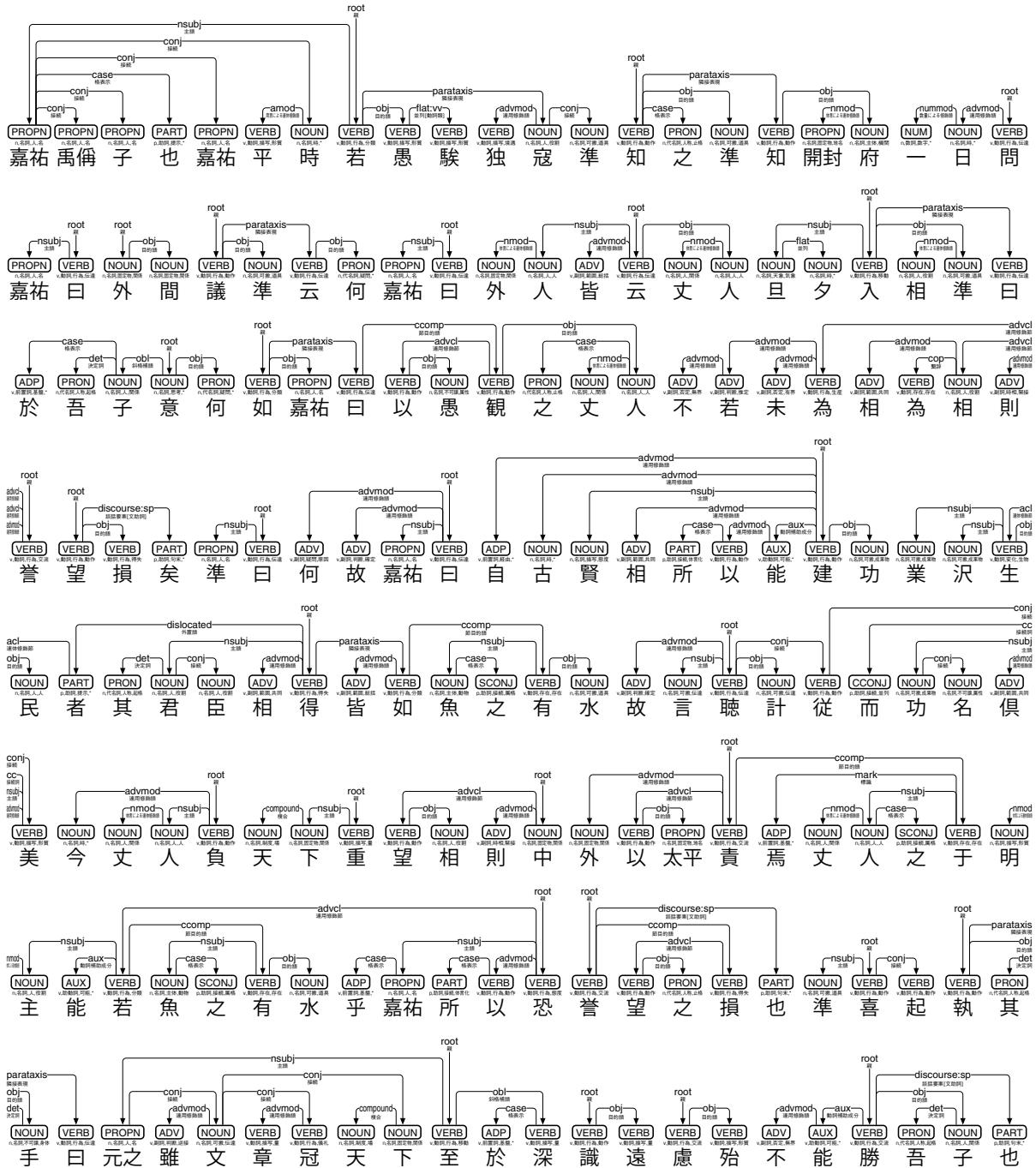


図 25: 手法②の処理結果(2018年)

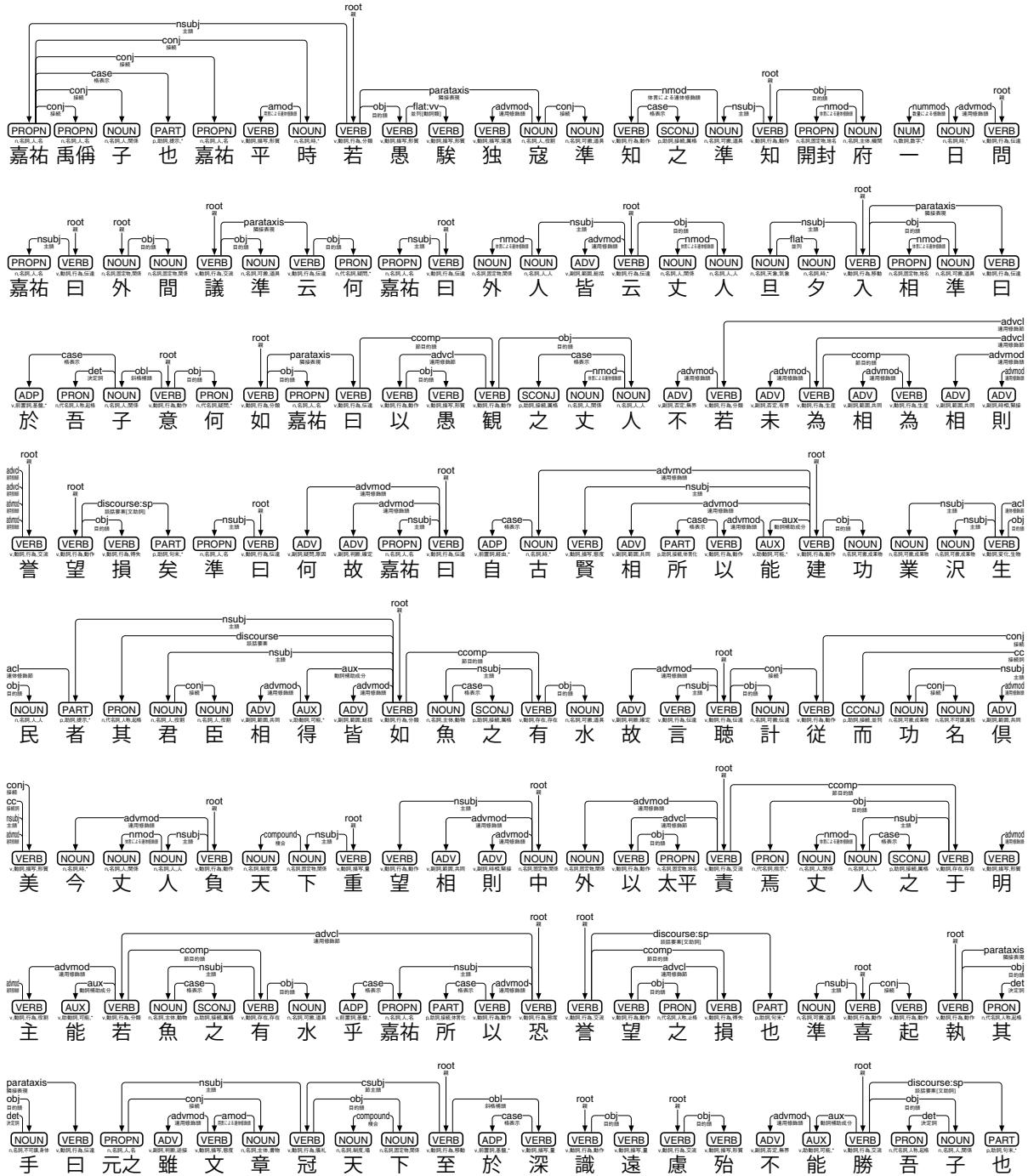


図 26: 手法③の処理結果 (2018 年)

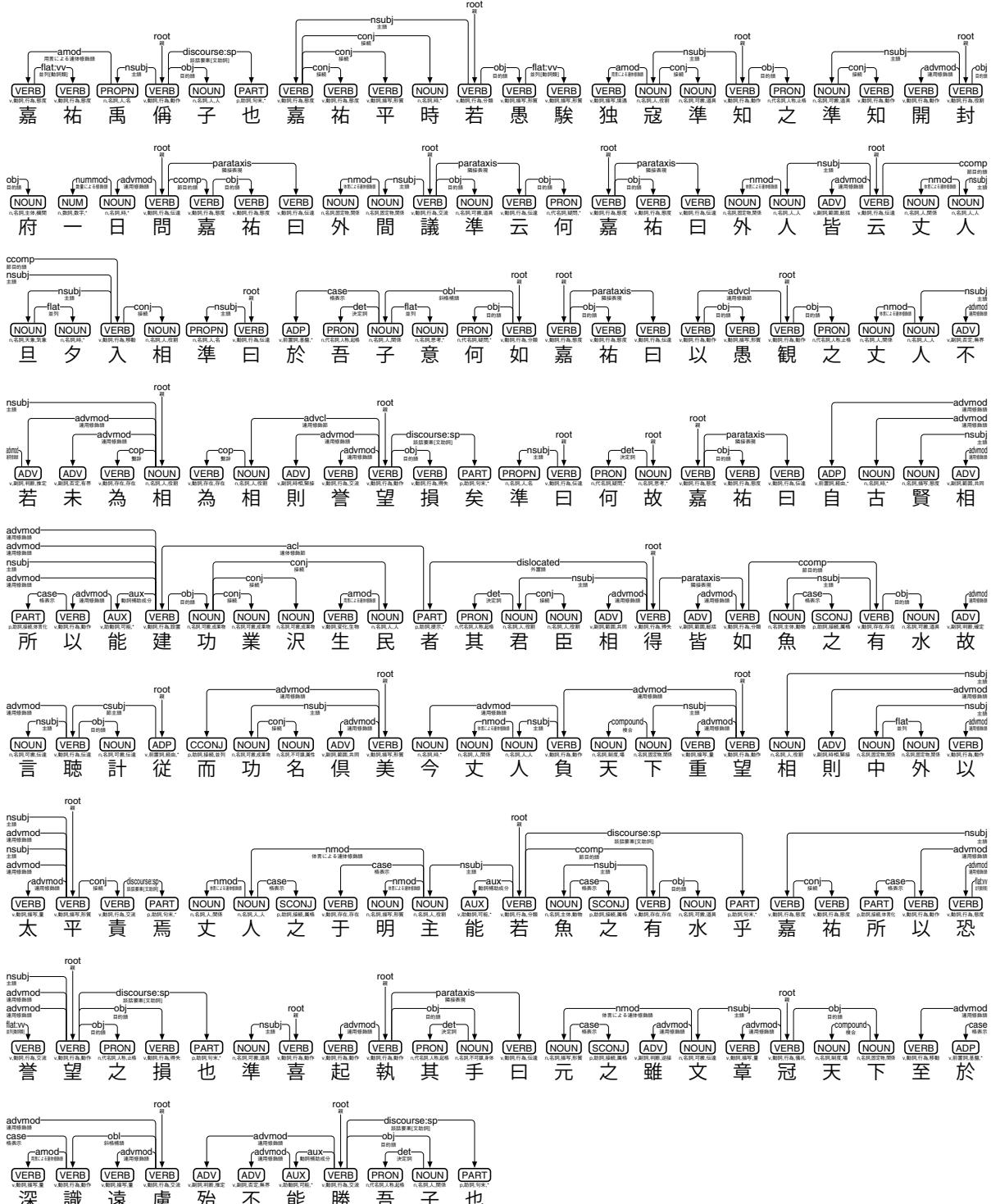


図 27: 手法①の処理結果(2018年)

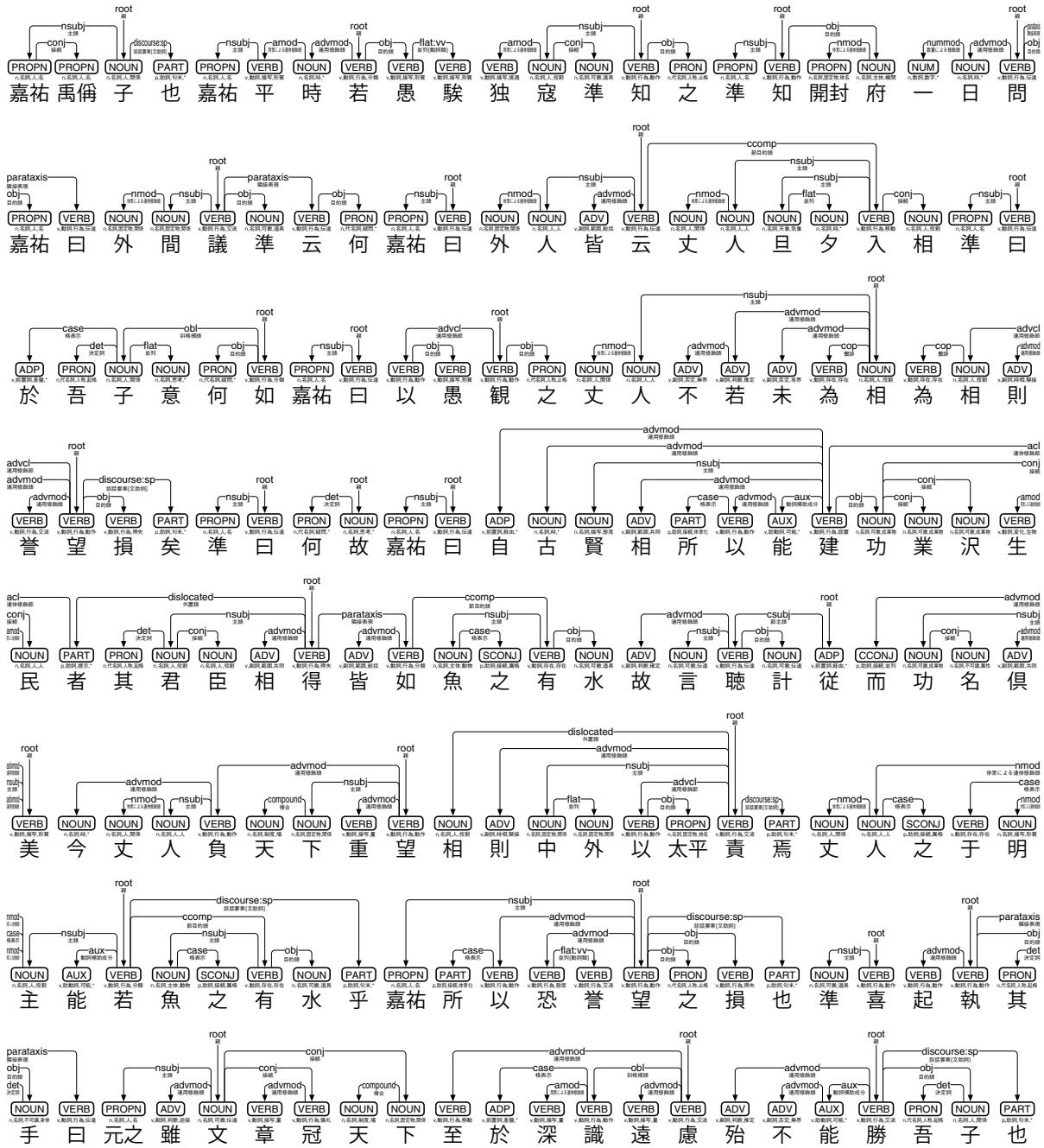


図 28: 手法②の処理結果(2018 年)

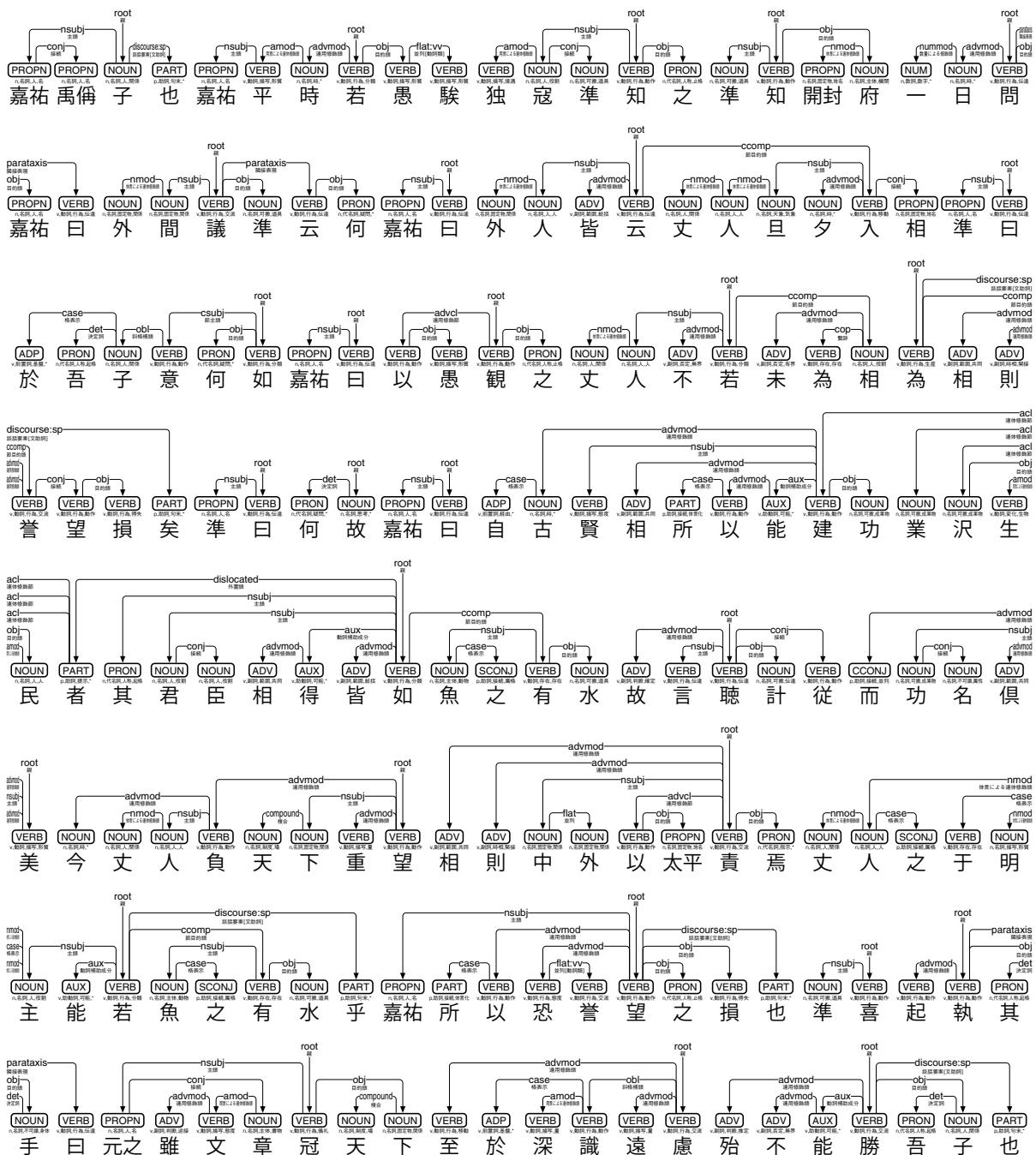


図 29: 手法③の処理結果(2018年)

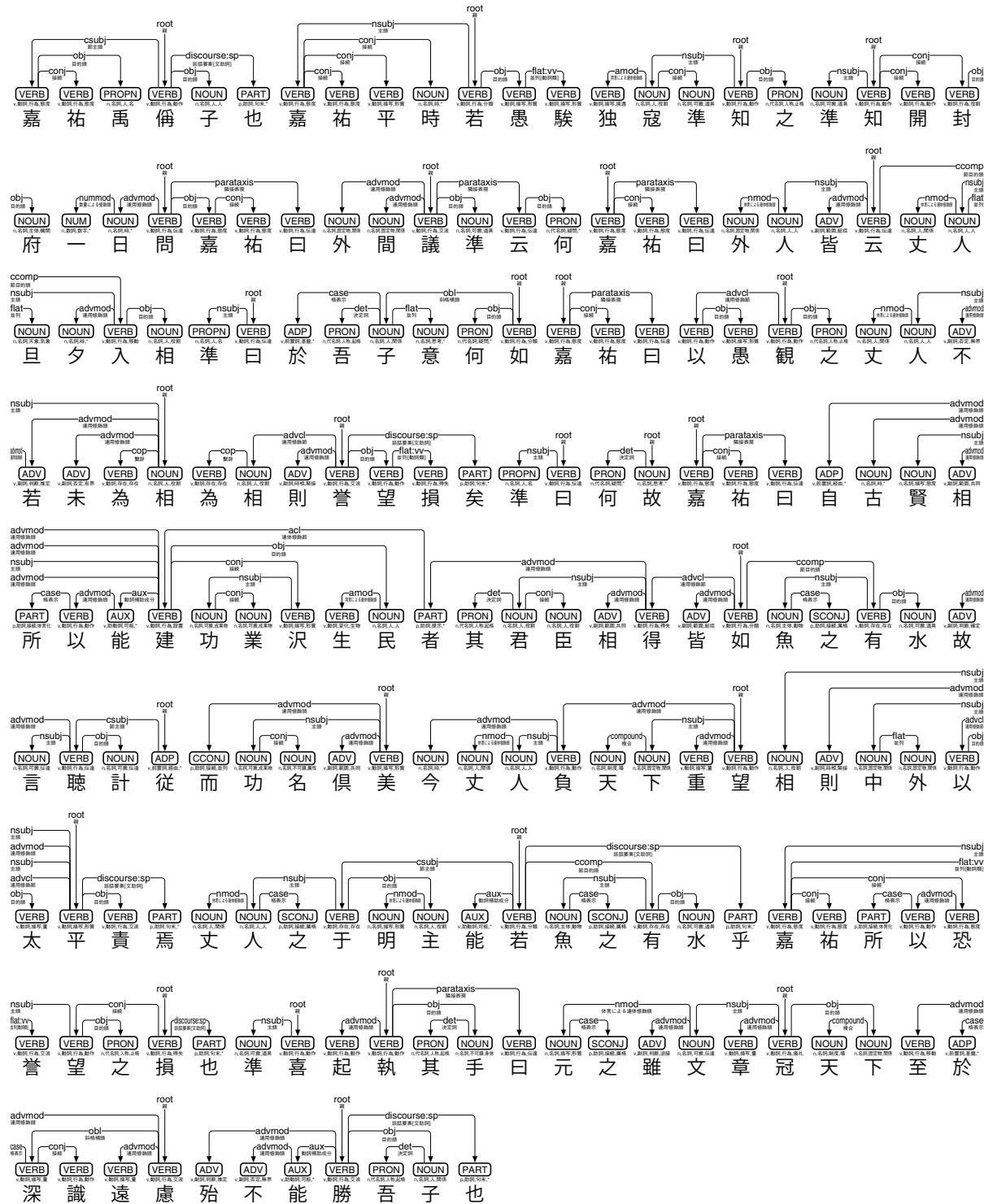


図 30: 手法①の処理結果(2018 年)

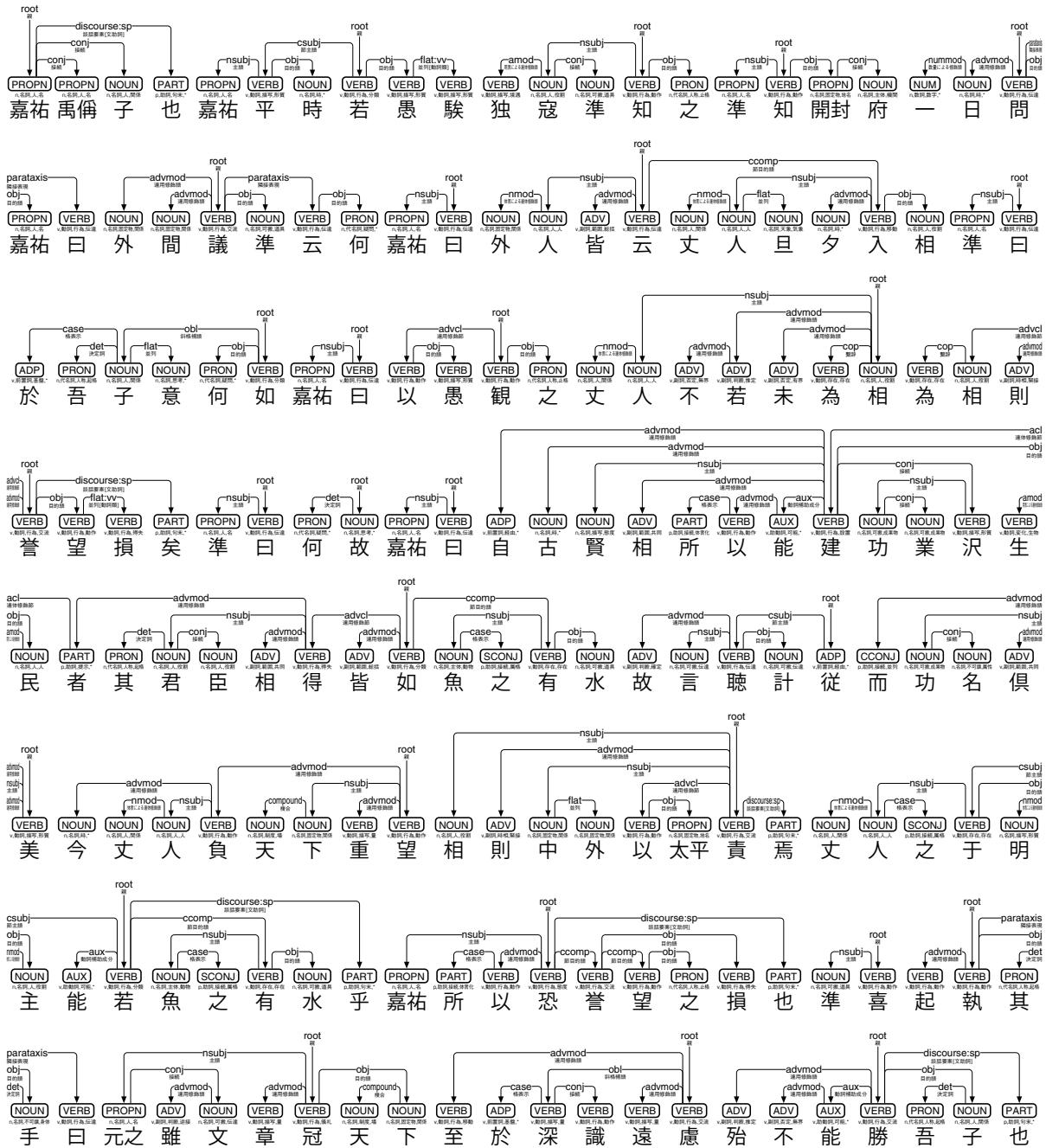


図 31: 手法②の処理結果(2018 年)

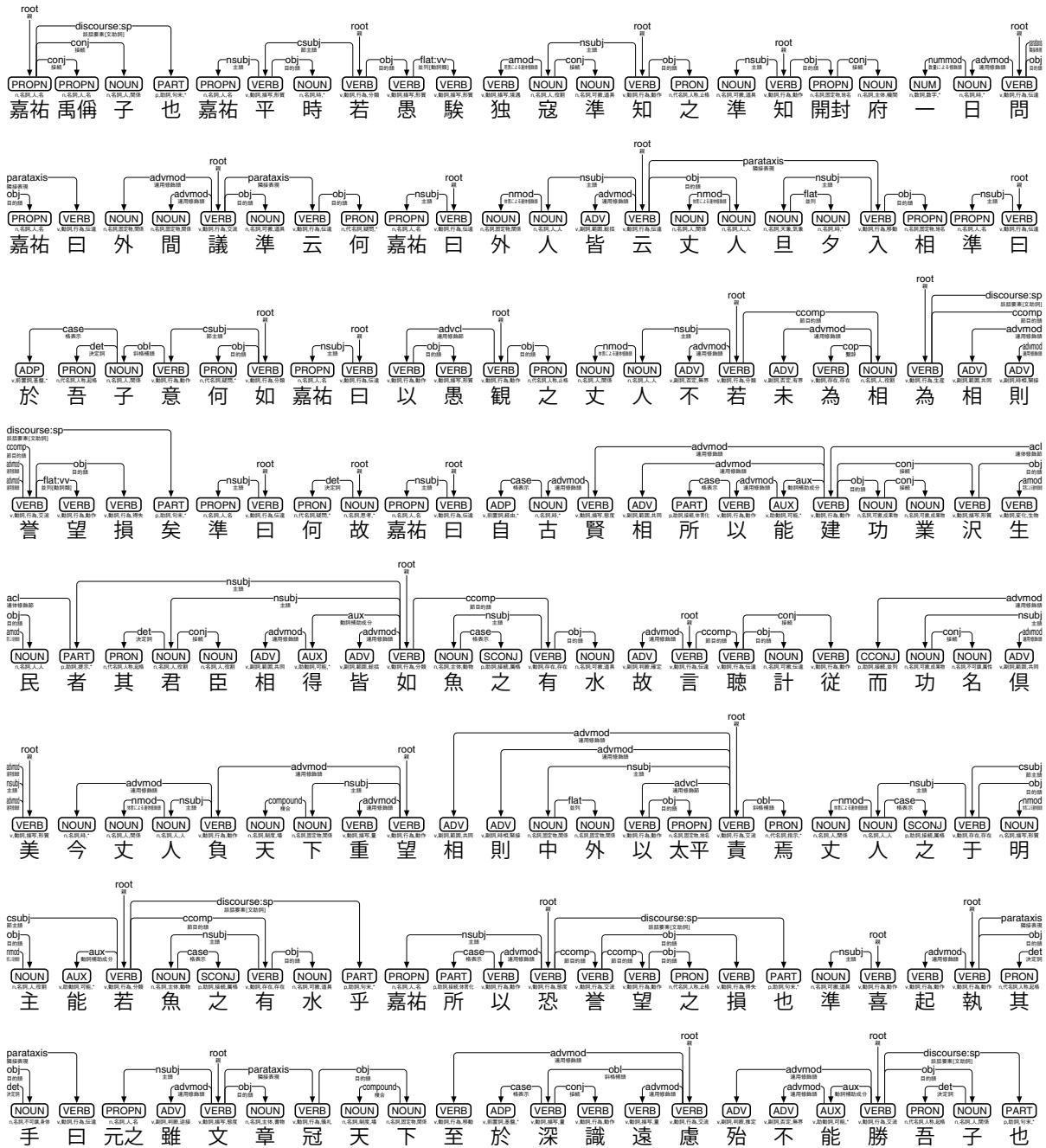


図 32: 手法③の処理結果(2018 年)

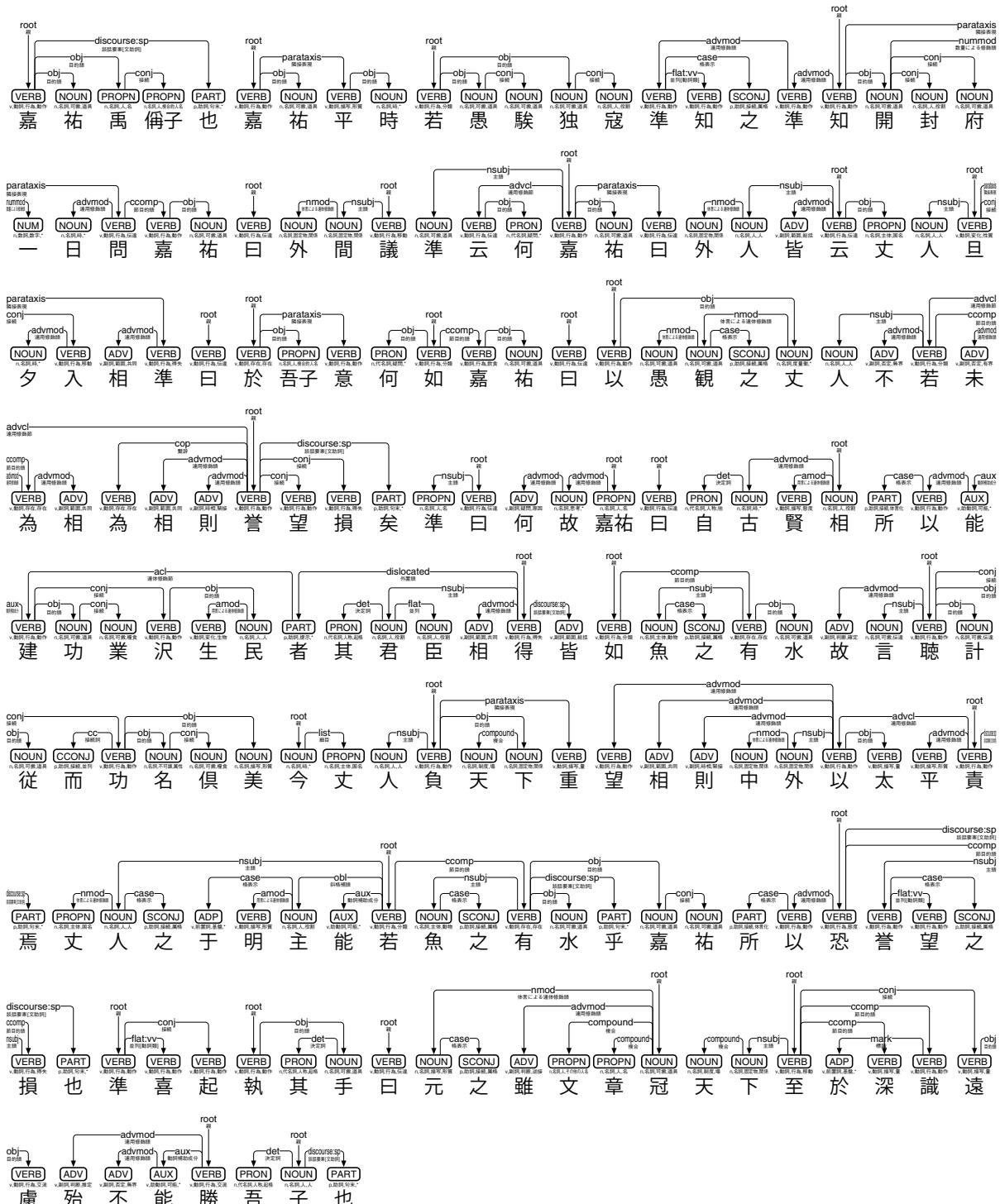


図 33: 手法Ⓐの処理結果(2018年)

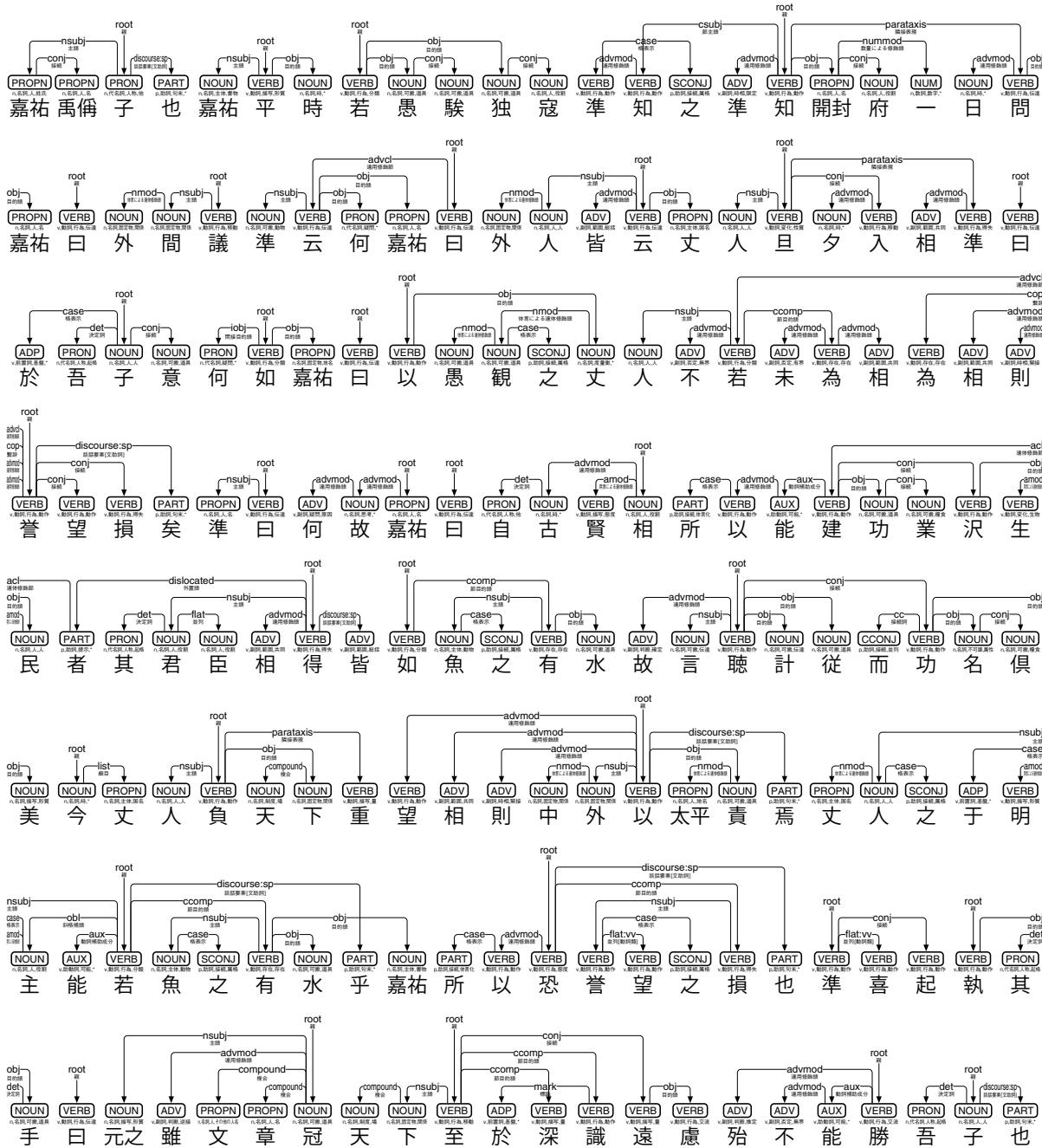


図 34: 手法②の処理結果(2018 年)

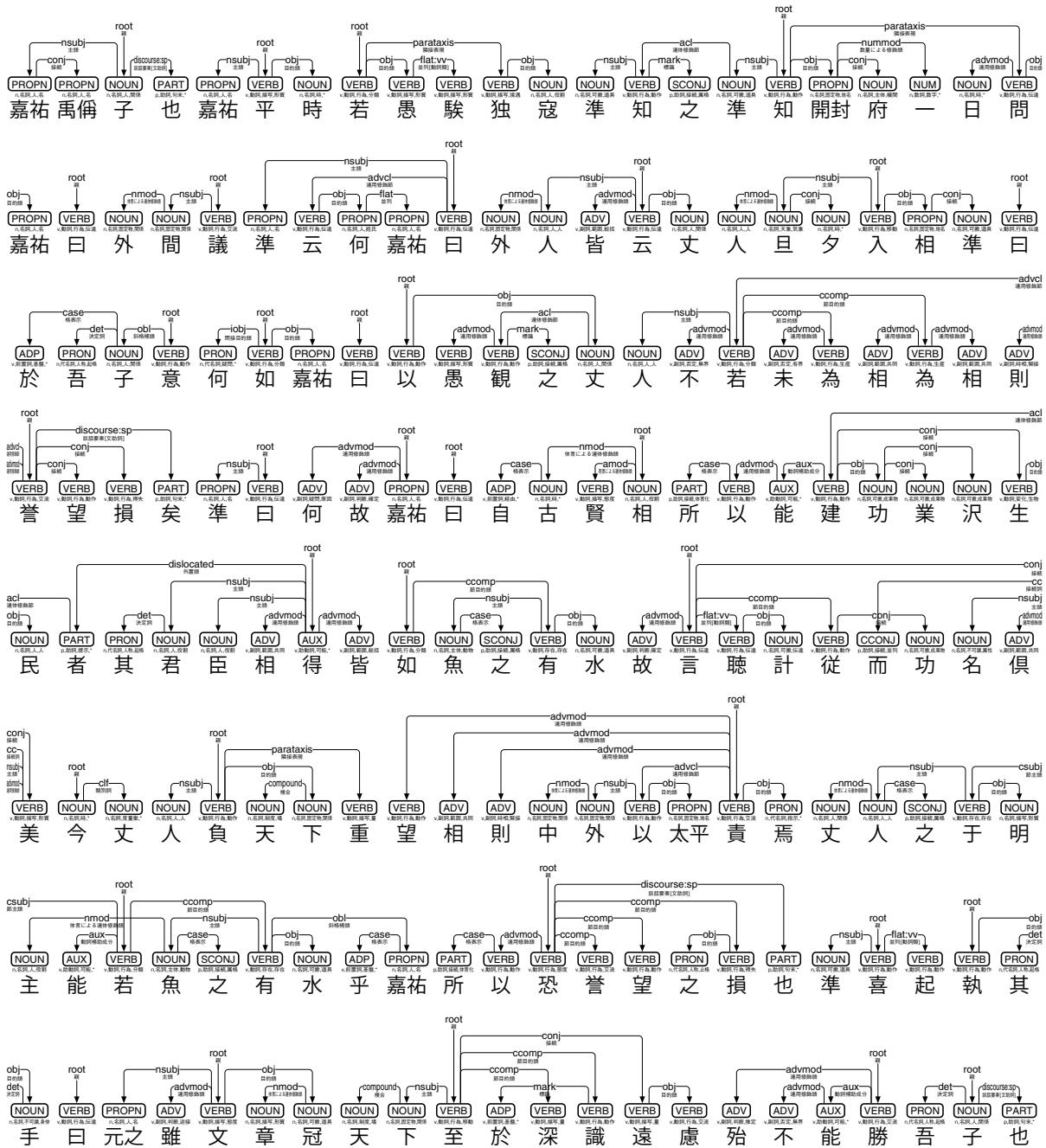


図 35: 手法④の処理結果(2018 年)

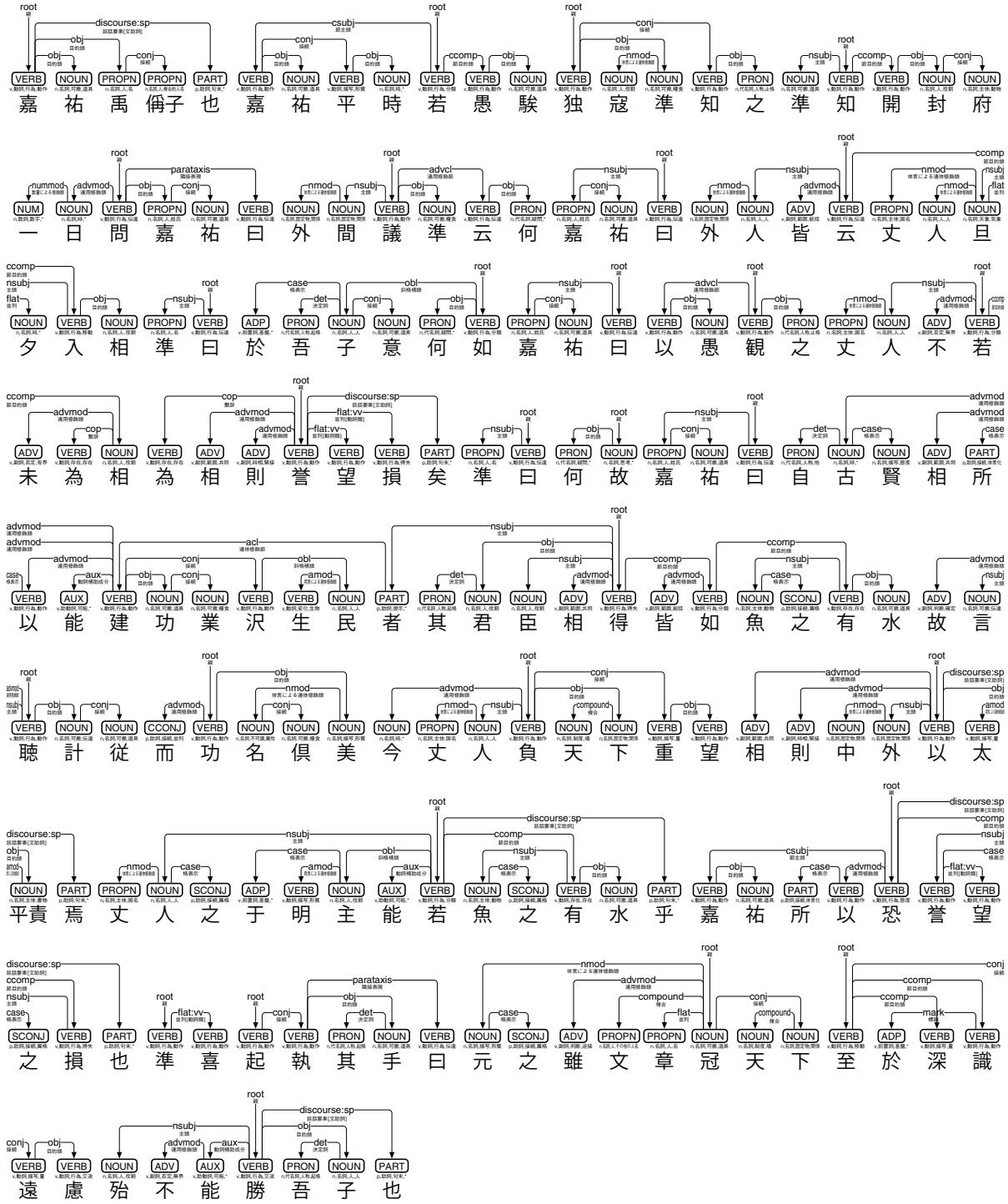


図 36: 手法Aの処理結果(2018年)

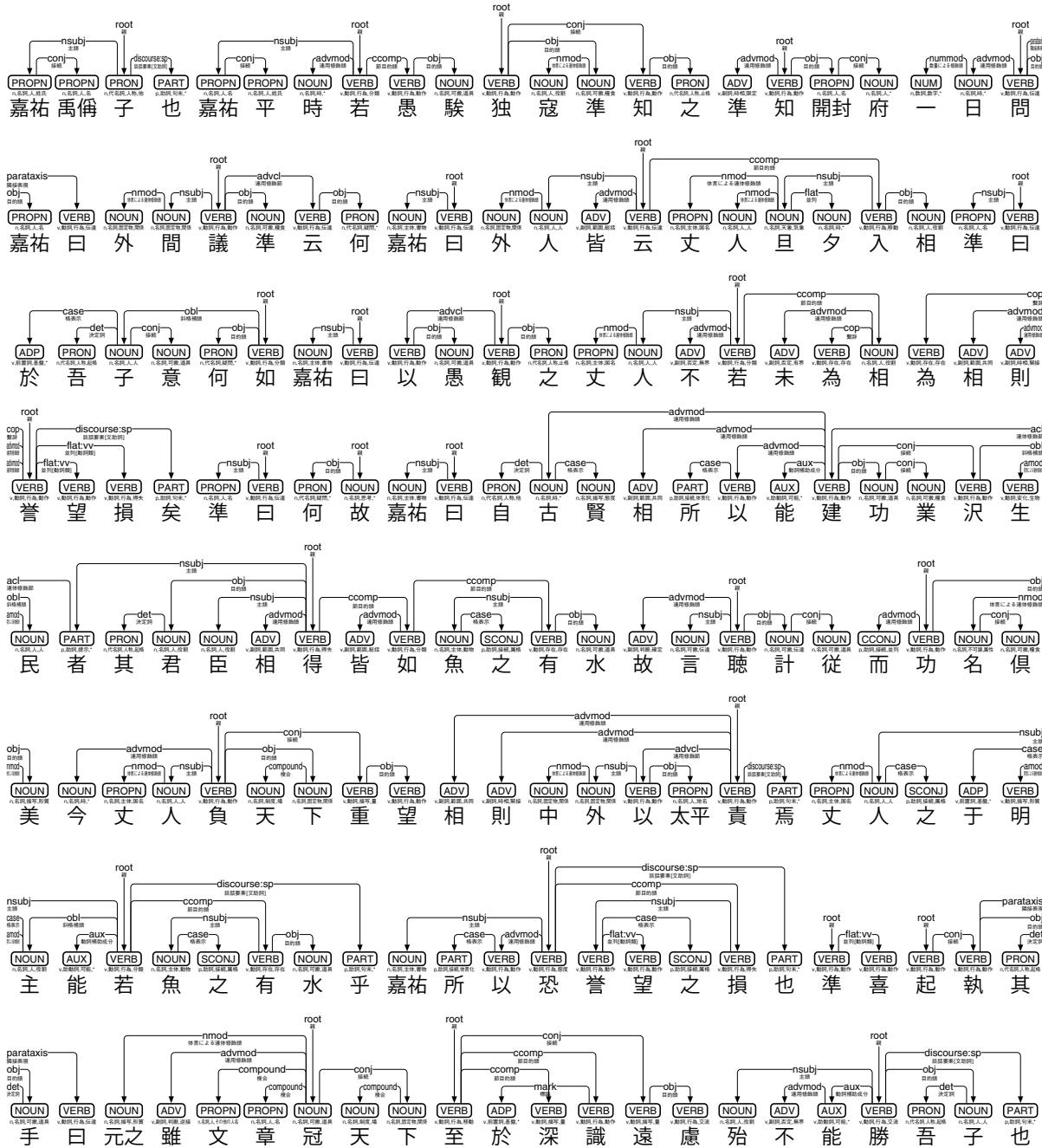


図 37: 手法Bの処理結果(2018年)

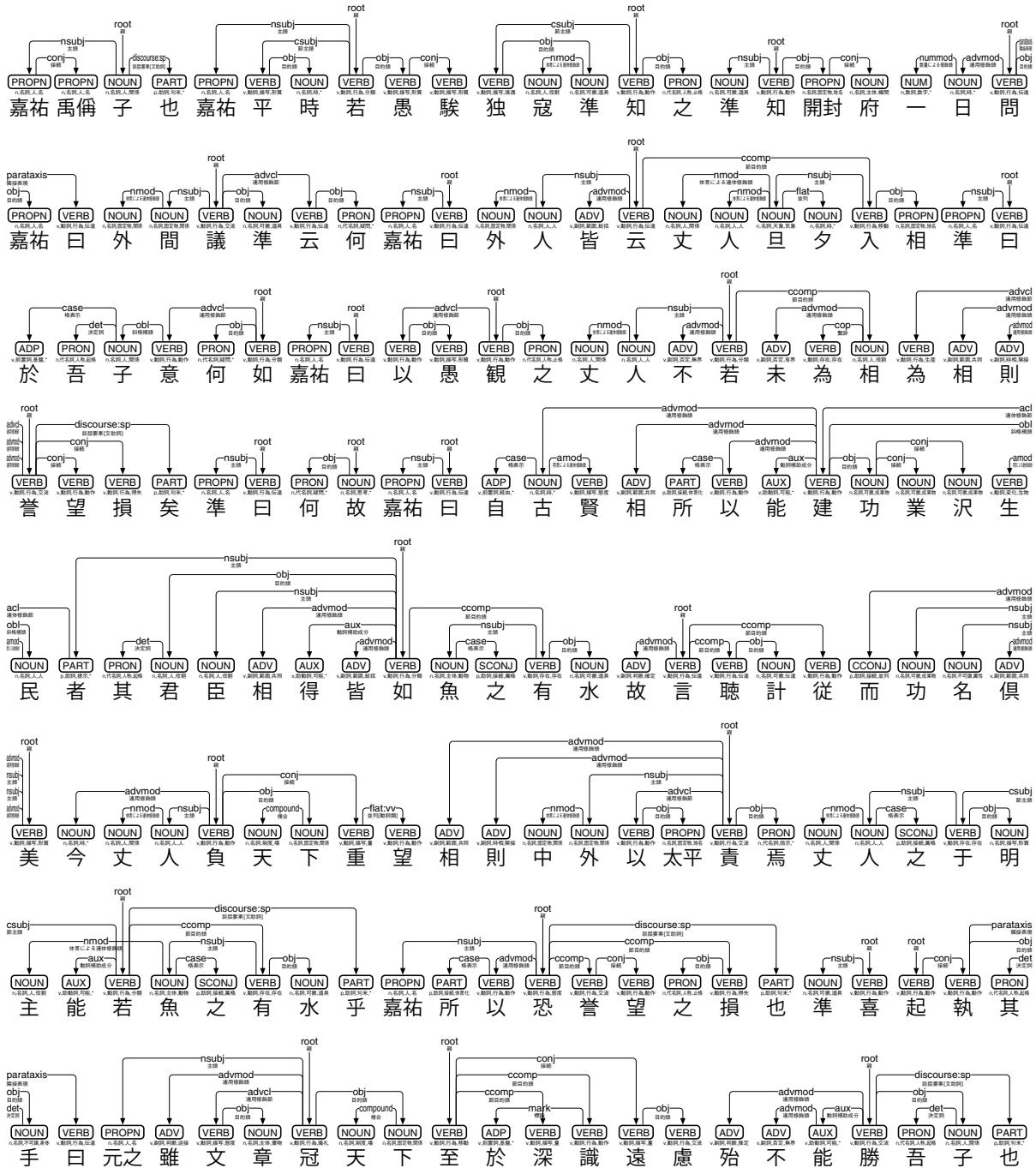


図 38: 手法Cの処理結果(2018年)

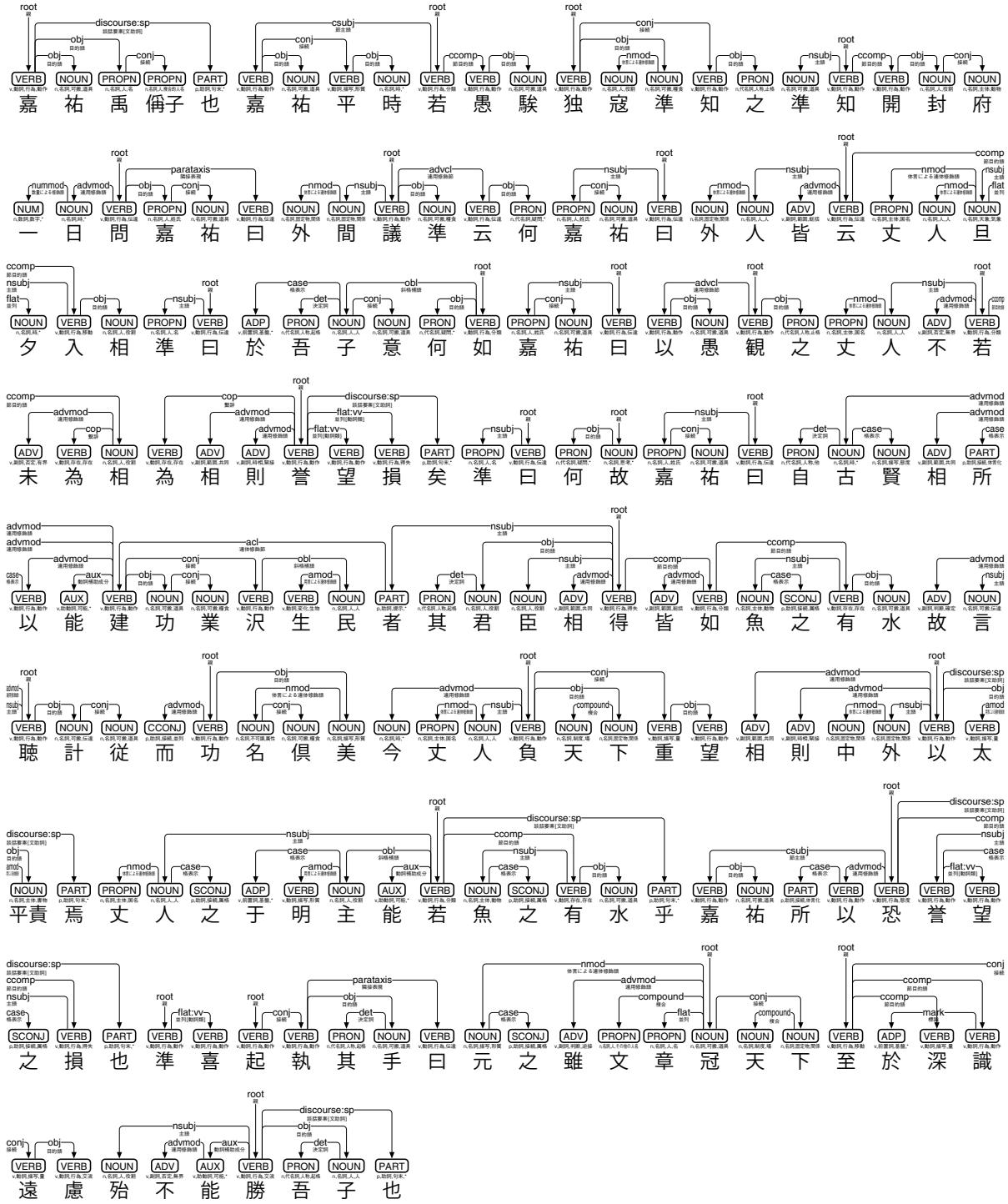


図 39: 手法Aの処理結果(2018年)

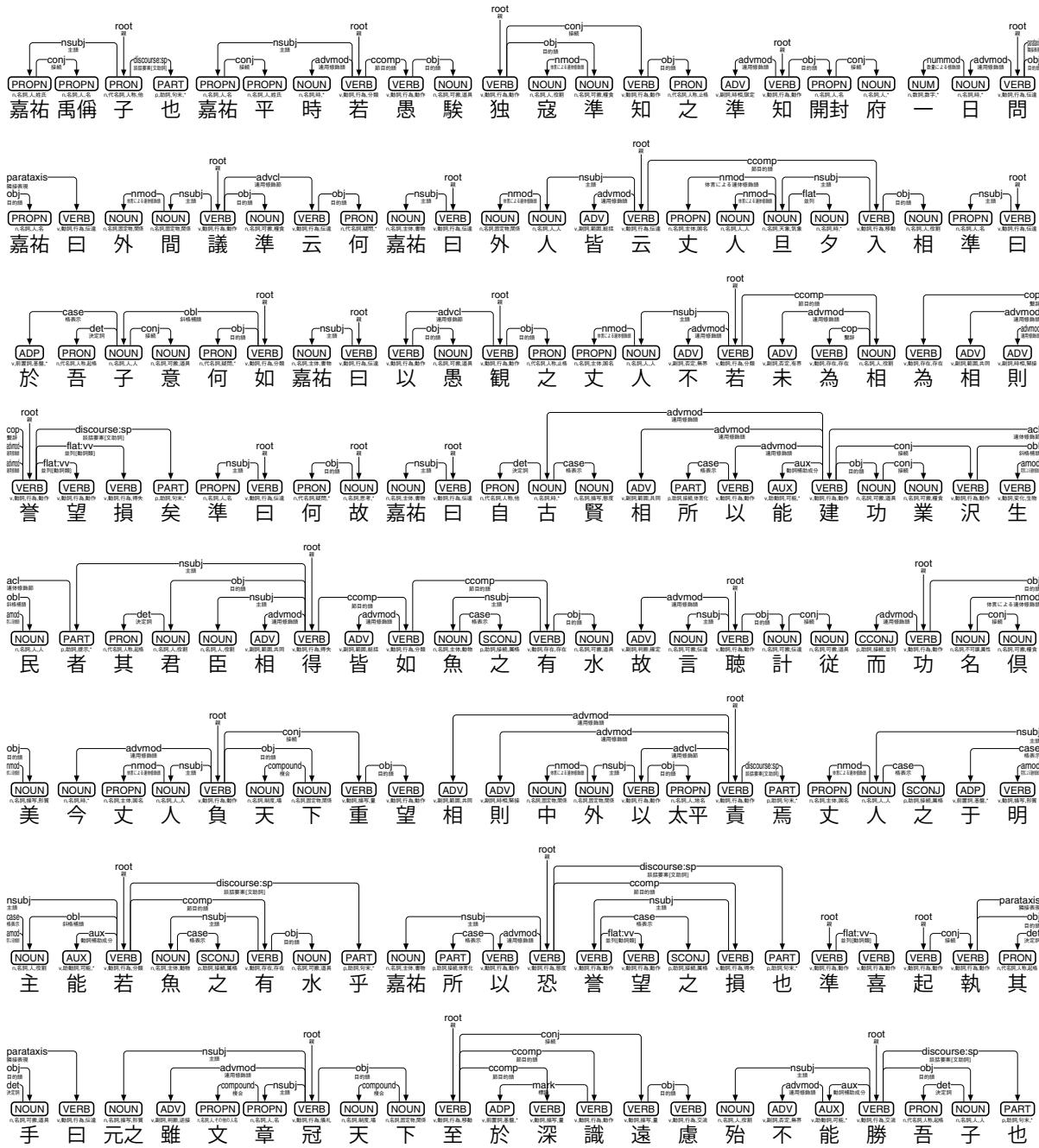


図 40: 手法Bの処理結果(2018年)

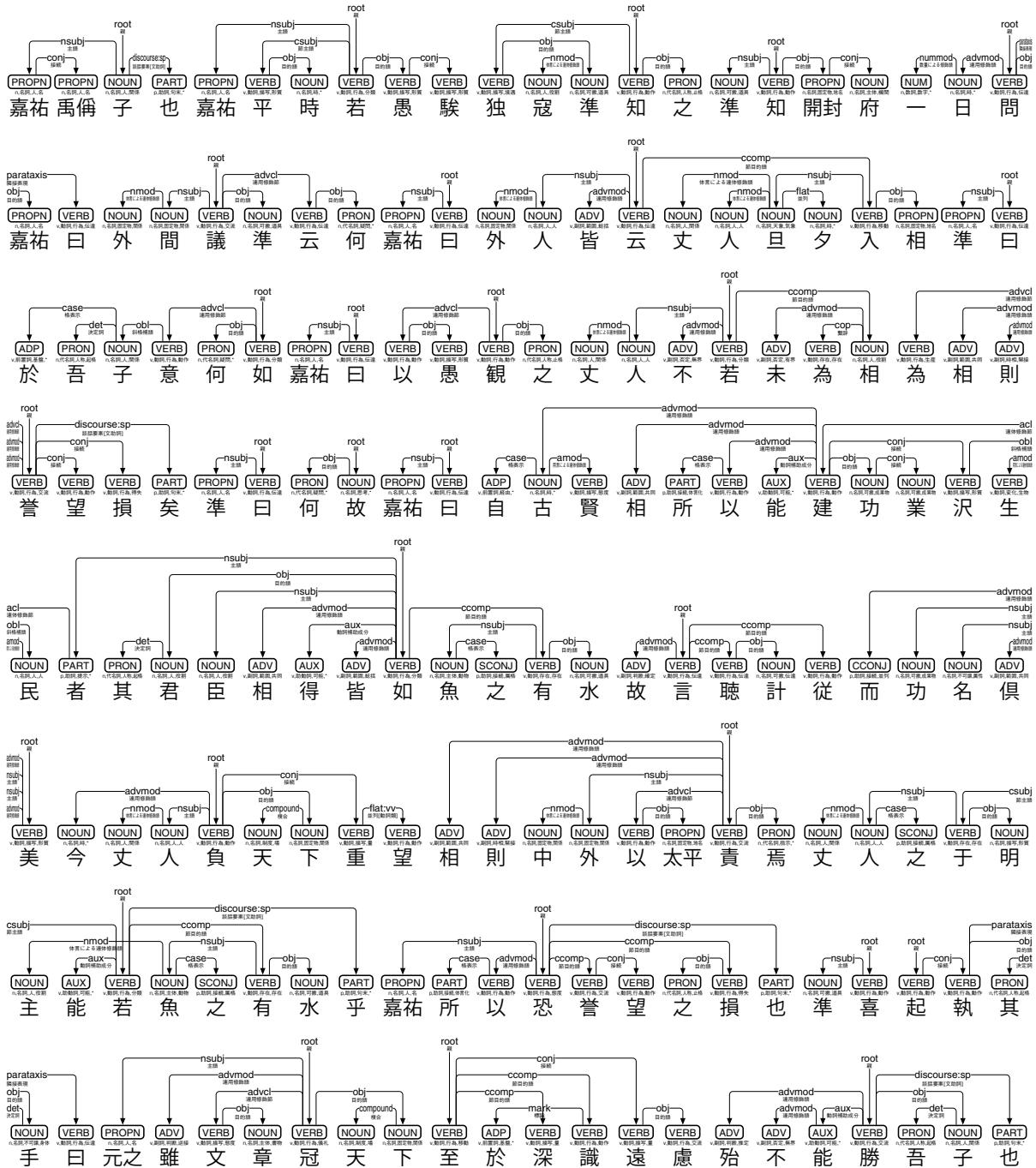


図 41: 手法Cの処理結果(2018年)

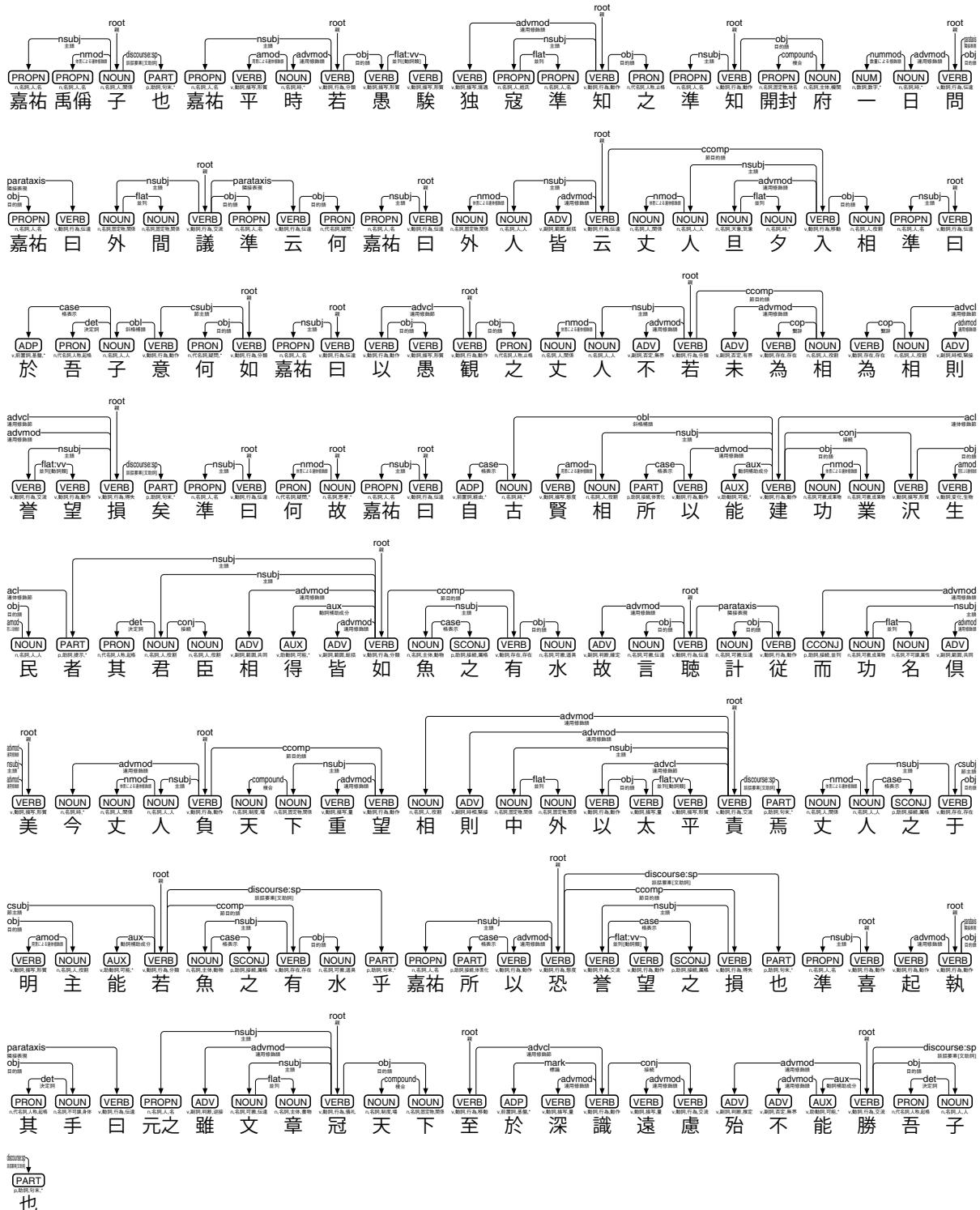


図 42: 手作業で作成した「正解」UD (2018 年)

A

聽ケバ^(注1) 雷ラバ 霆テイヲ 於百里之外ニ者、如ク 鼓スルガ^(注2) 盆ヲ 望メバ 江河ヲ 於千里之間ニ者、

如キハ 繫マトブガ 帶ヲ 以テ 其ノ 相ヒ 去ル 之ヲ 遠キヲ 也。故ニ 居リテ 千載ヲ 下ニ 而ニ 求ムルニ 之ヲ 千

載シテ 之ヲ 上ニ 以ニ 相ヒ 去ル 之ヲ 遠キヲ 而レバ 不レバ 知ラ 有ルヲ 其ノ 變チ 則シニ 猶ホ 刻ミシハリ 舟ヲ 求ムルガ 剣ヲ 今ニ 之ヲ

所ハ 求ムル 非ザルモ 往スル 者ヲ 所ハ 失フ 而レバ 謂おもへリ 下ノ 其ノ 刻ミシハリ 在リ 此ニ 是レ 所ニ 徒ヨリテ 墜オツル 也上。豈ニ 不レバ 感ヒナラ 乎ヤ。

C

今夫江戸者、世之所ハ 称スル 名都大邑だい(注4) 冠蓋いふ(注5) 之所ハ 集マル(2) 舟車ヲ 之ヲ

所ニシテ 湿あつ 実タル 為タル 天下之大都會タリ 也。而レドモ 其ノ 地ノ 為スル 名ヲ 訪フ 之ヲ 古ニ 未タリ

之ヲ 聞フ。豈ニ 非ズ 古今相ヒ 去ルコト(1) 日ひびニ 遠ク 而レバ 事物之變モ 亦タリ 在リ 于ノ 其ノ 間ヲ 耶や(ア) 蓋シ

知ル 後ニ 之ヲ 於カルモ 今ニ 世ノ 相ヒ 去ルコト(1) 愈スル 遠ク 事ノ 之ヲ 變コト 愈スル 多ク 求ムルモ 其ノ 所ヲ 欲スル 聞カント

而レバ 不ルコト 可カラ 得タ 亦猶 今ニ 之ヲ 於カルガ 古ニ 也。

B

D

吾ひそかニ 有リ 感ズル 焉。 遺ロ 聞ム 之ヲ 書ム 所ニ 由ヨリテ 作ル 也。

図 43: 大学入試センター試験『国語』(2017年1月14日)第4問本文

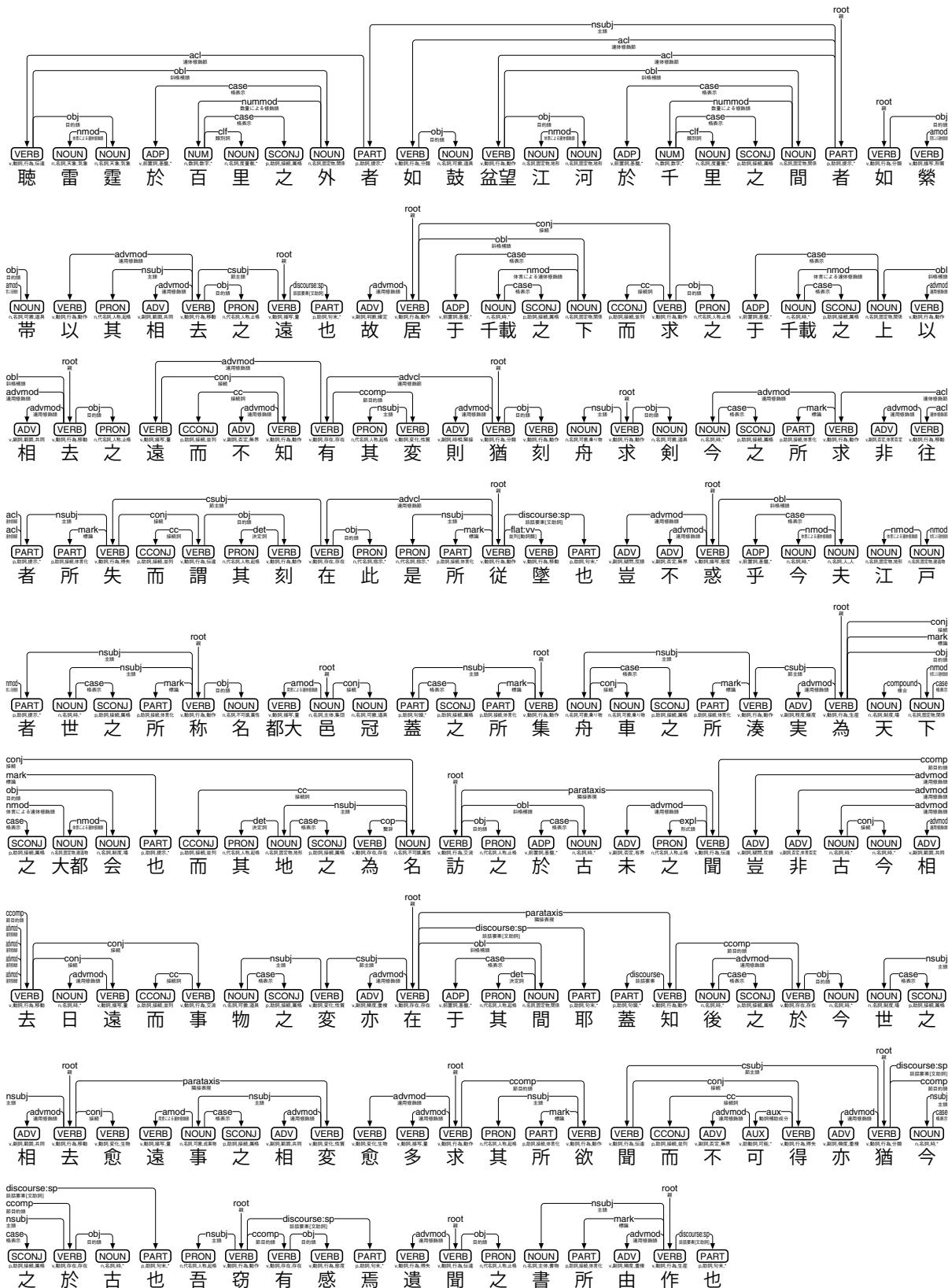


図 44: 手法①の処理結果(2017 年)

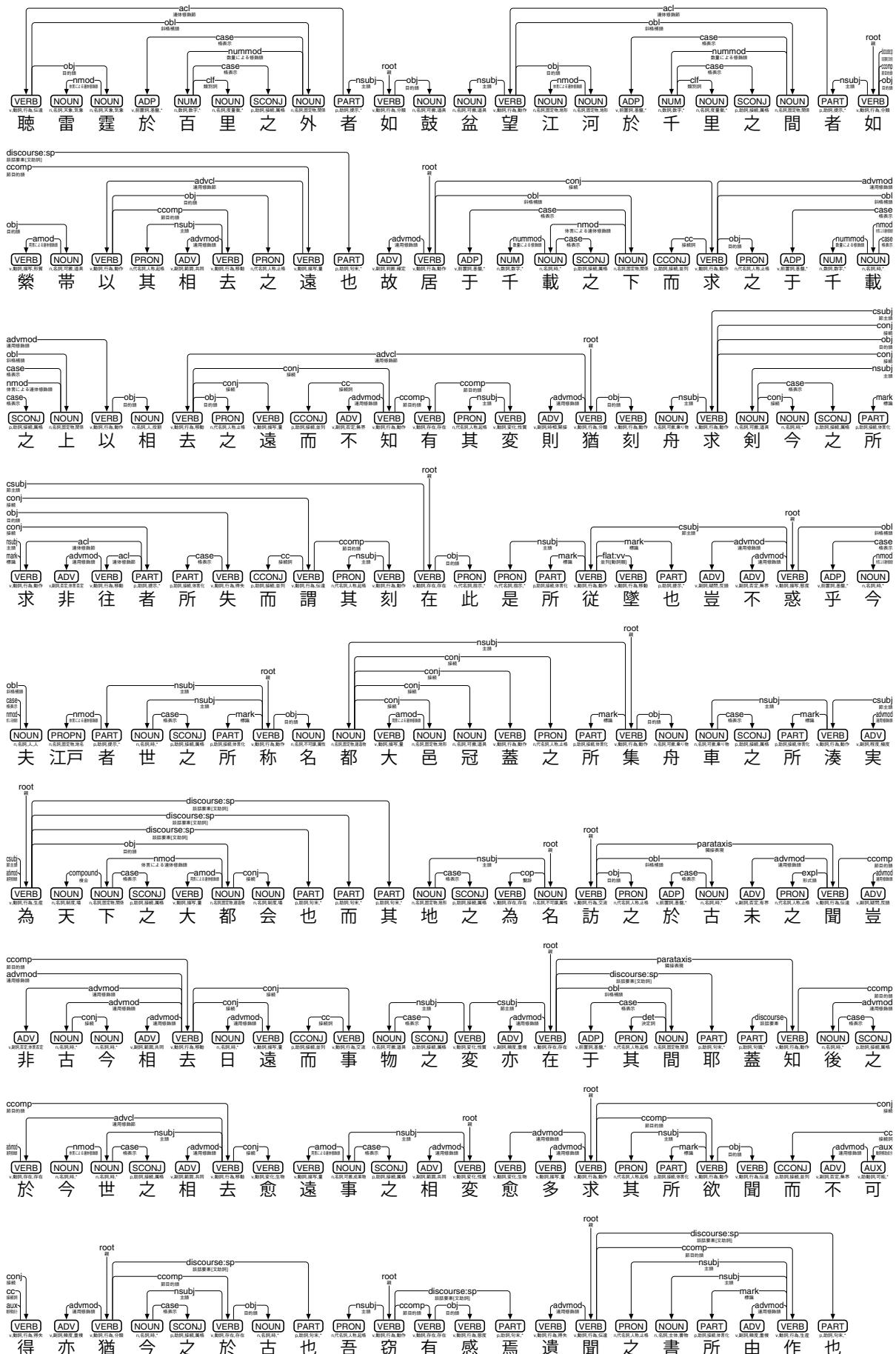


図 45: 手法②の処理結果(2017年)

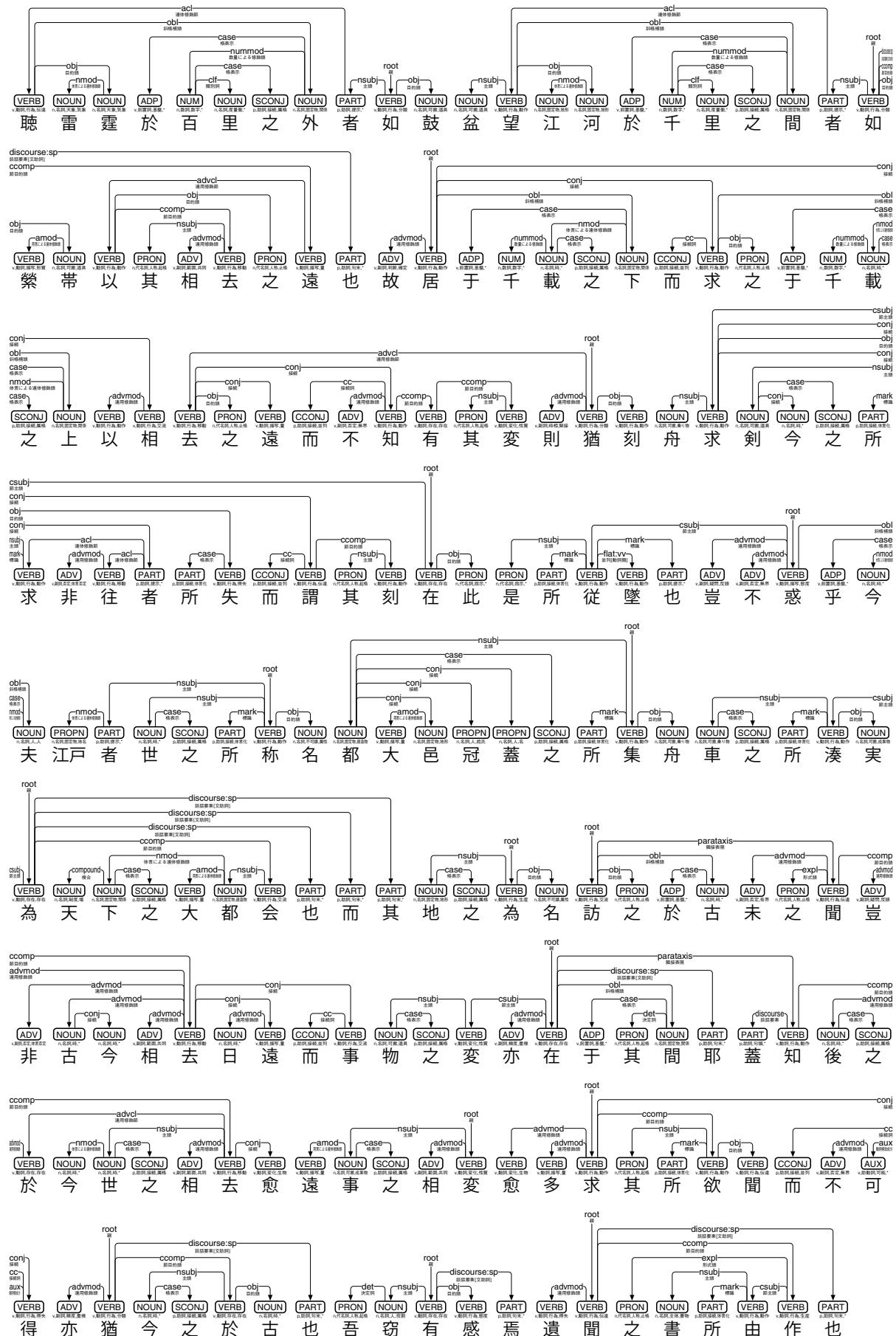


図 46: 手法③の処理結果 (2017 年)

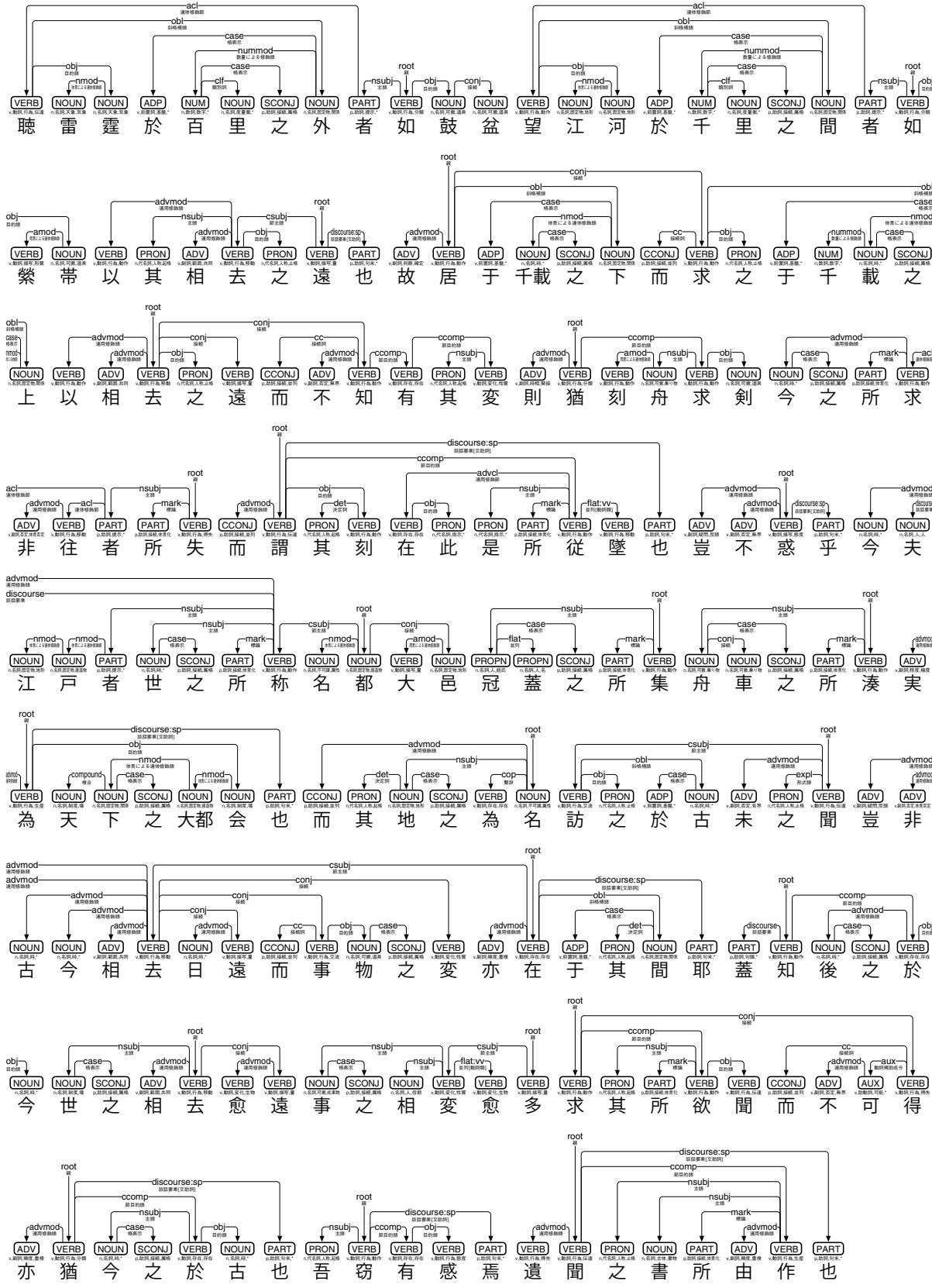


図 47: 手法①の処理結果(2017年)

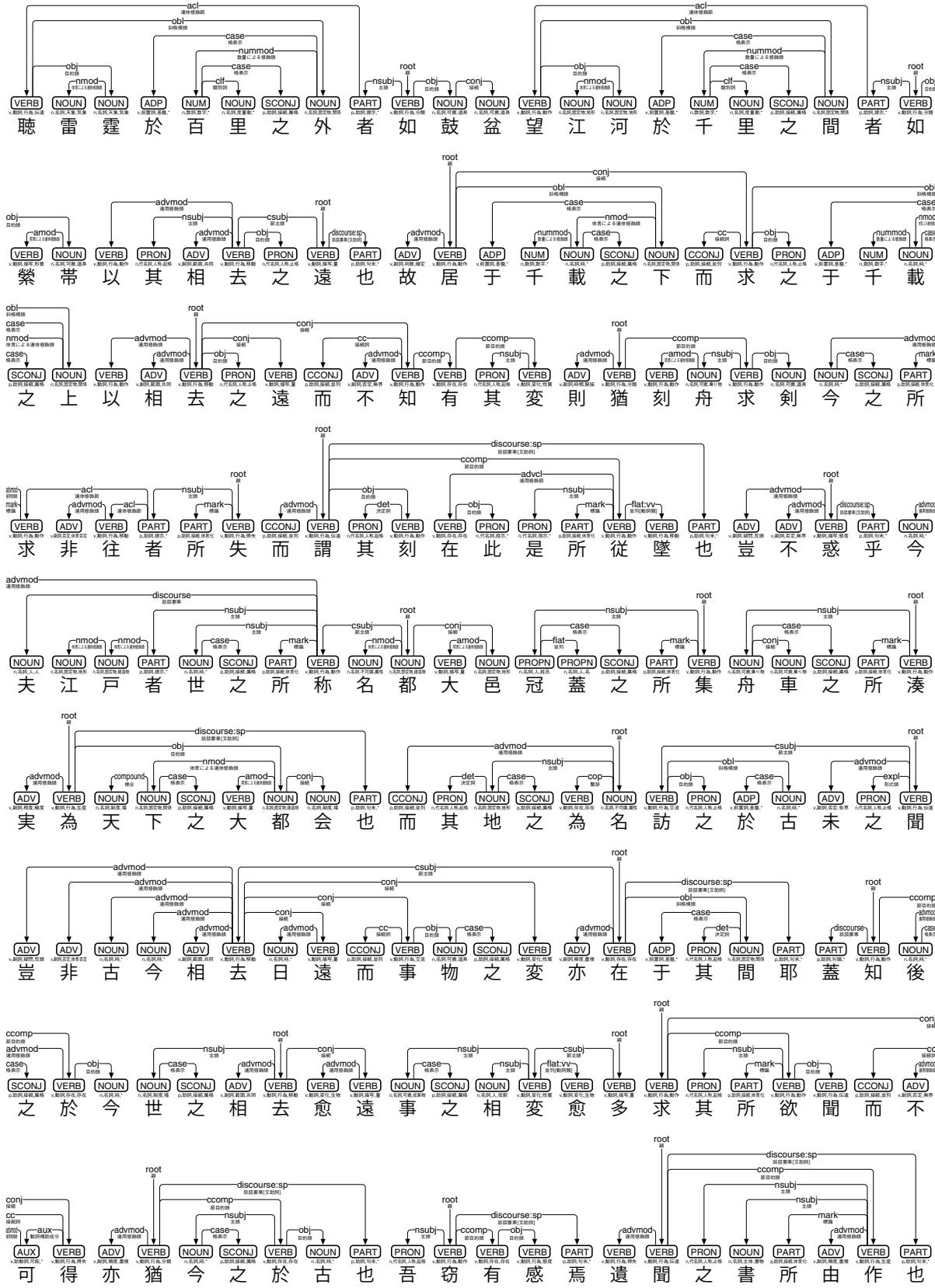


図 48: 手法②の処理結果(2017年)

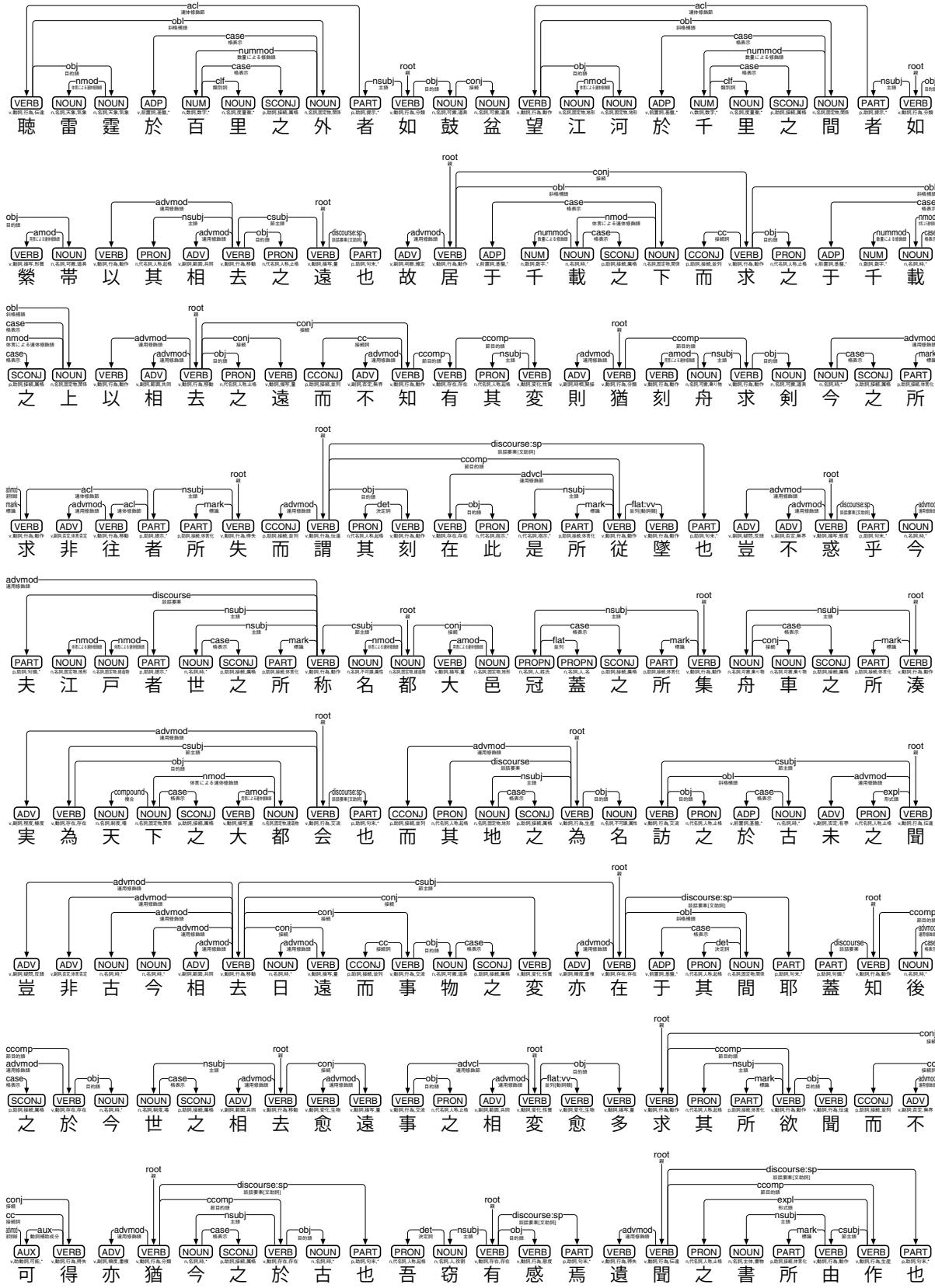


図 49: 手法③の処理結果(2017年)

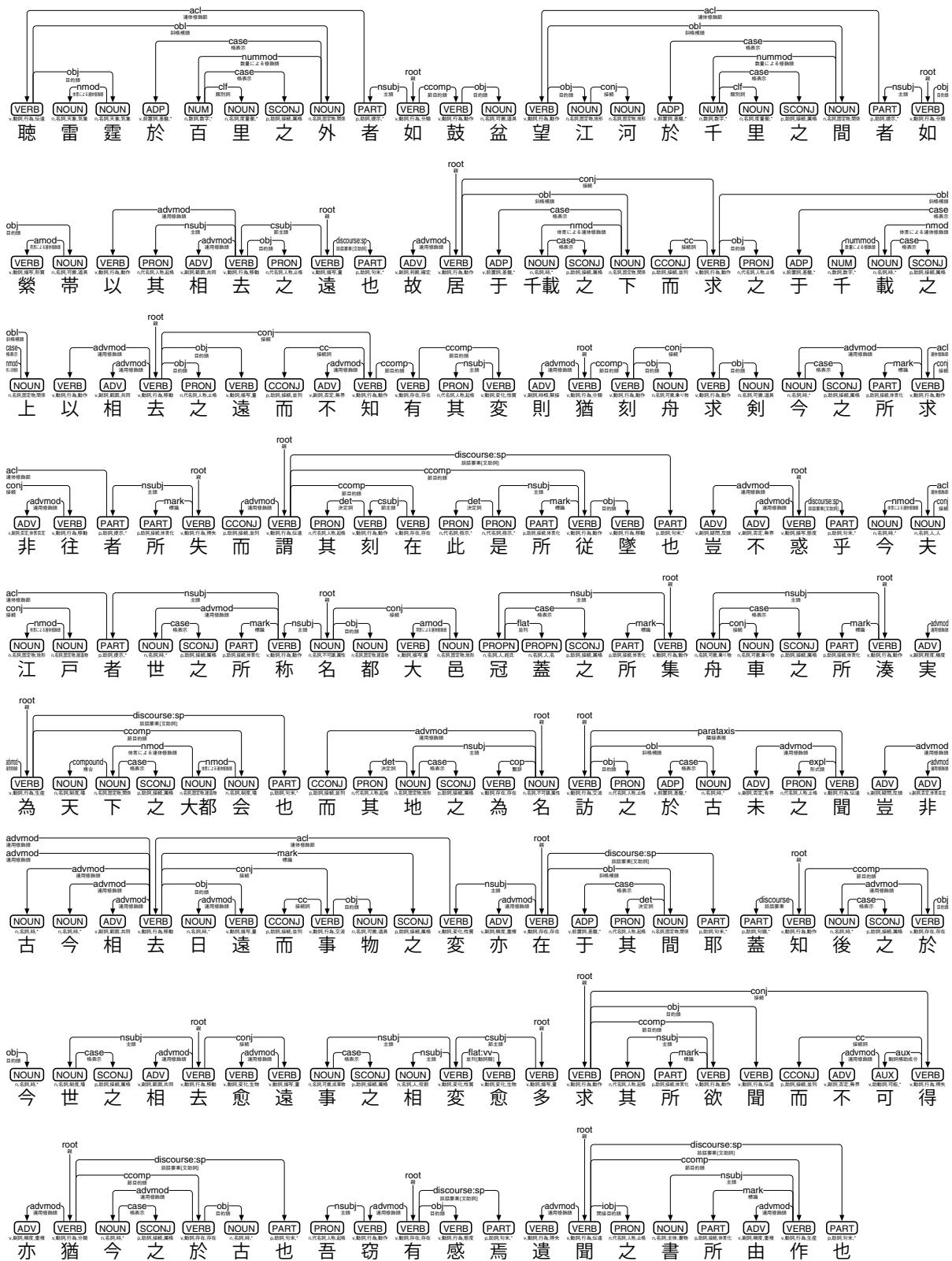


図 50: 手法①の処理結果(2017 年)

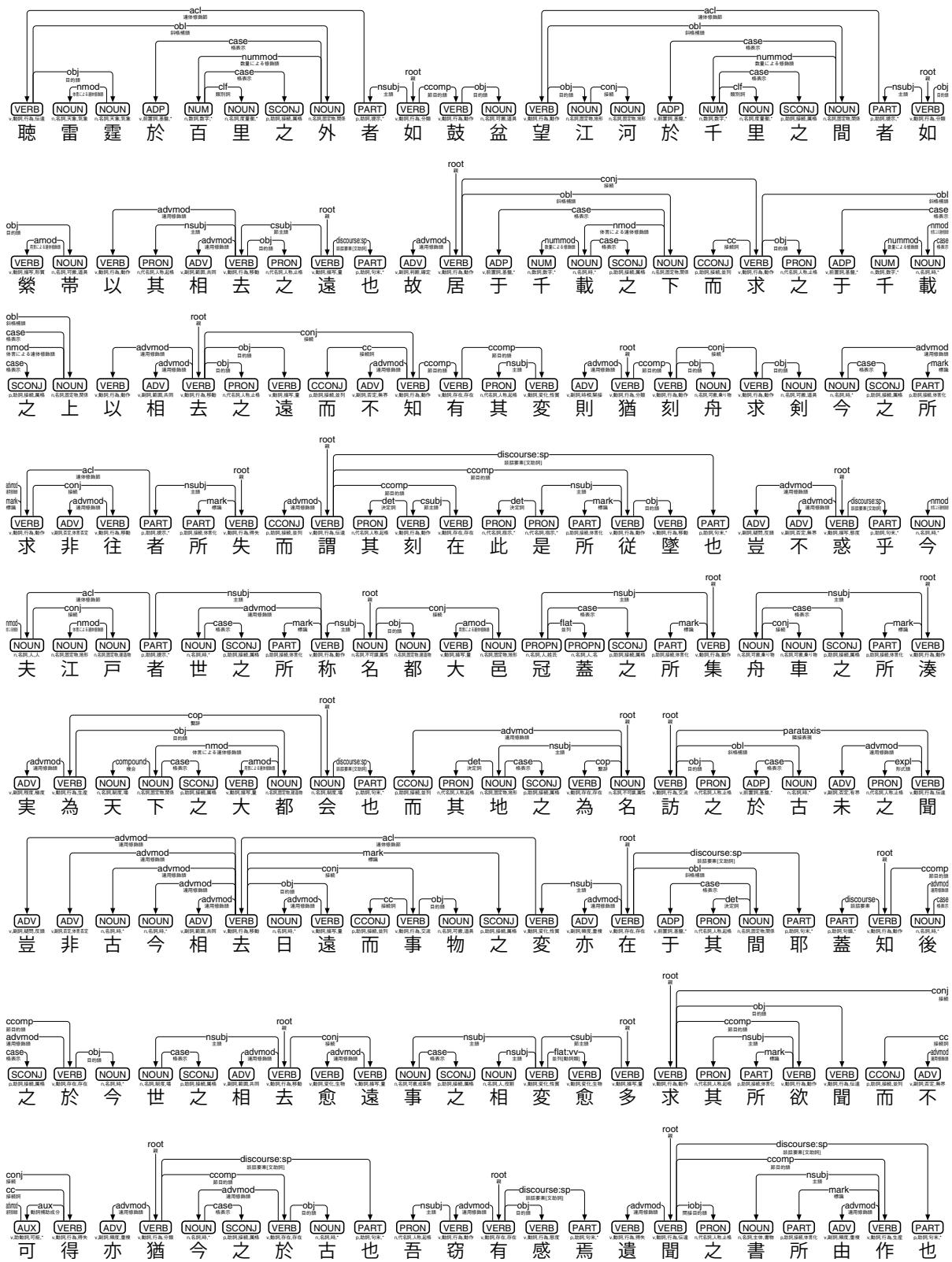


図 51: 手法②の処理結果(2017 年)

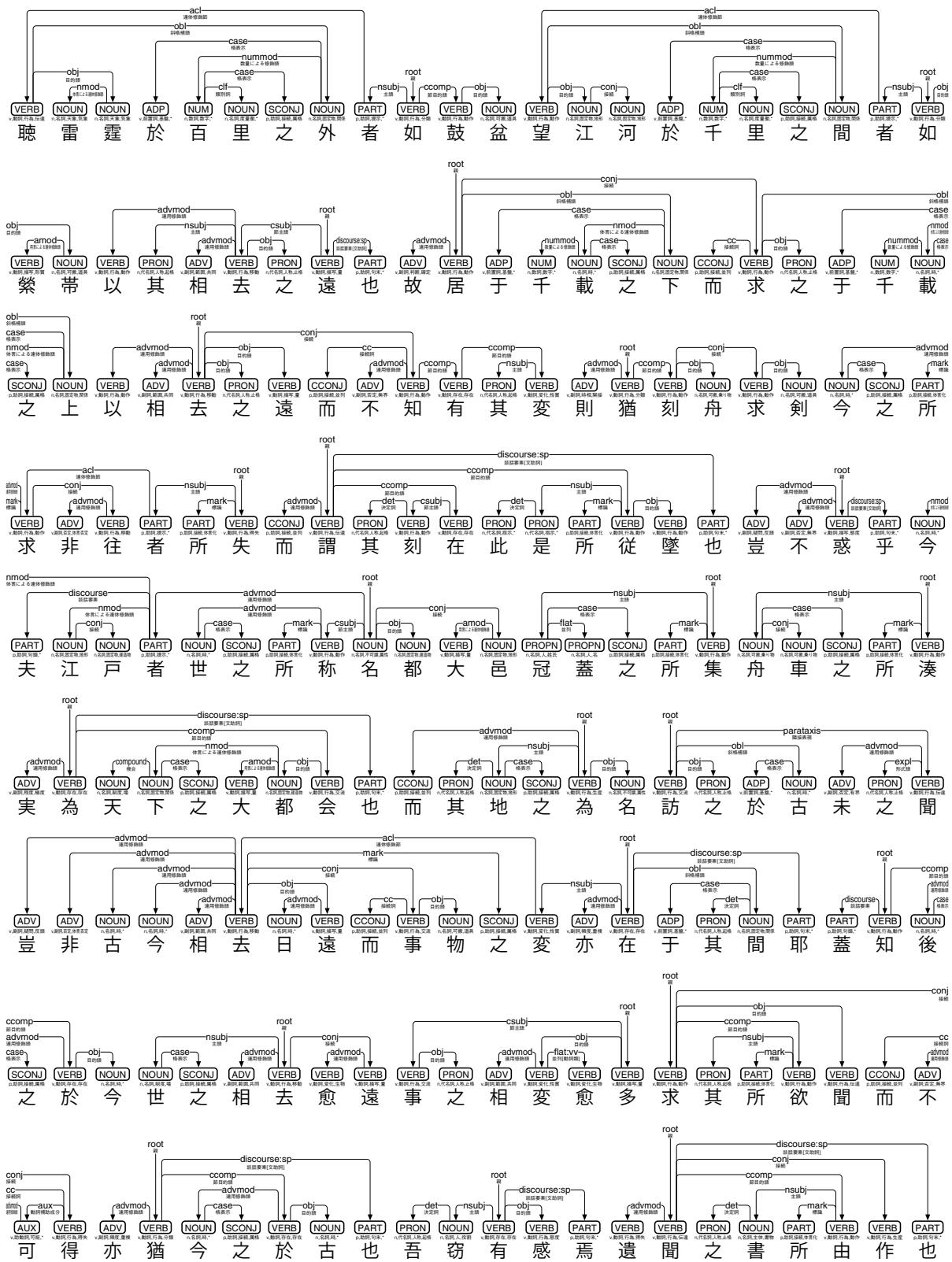


図 52: 手法③の処理結果(2017 年)

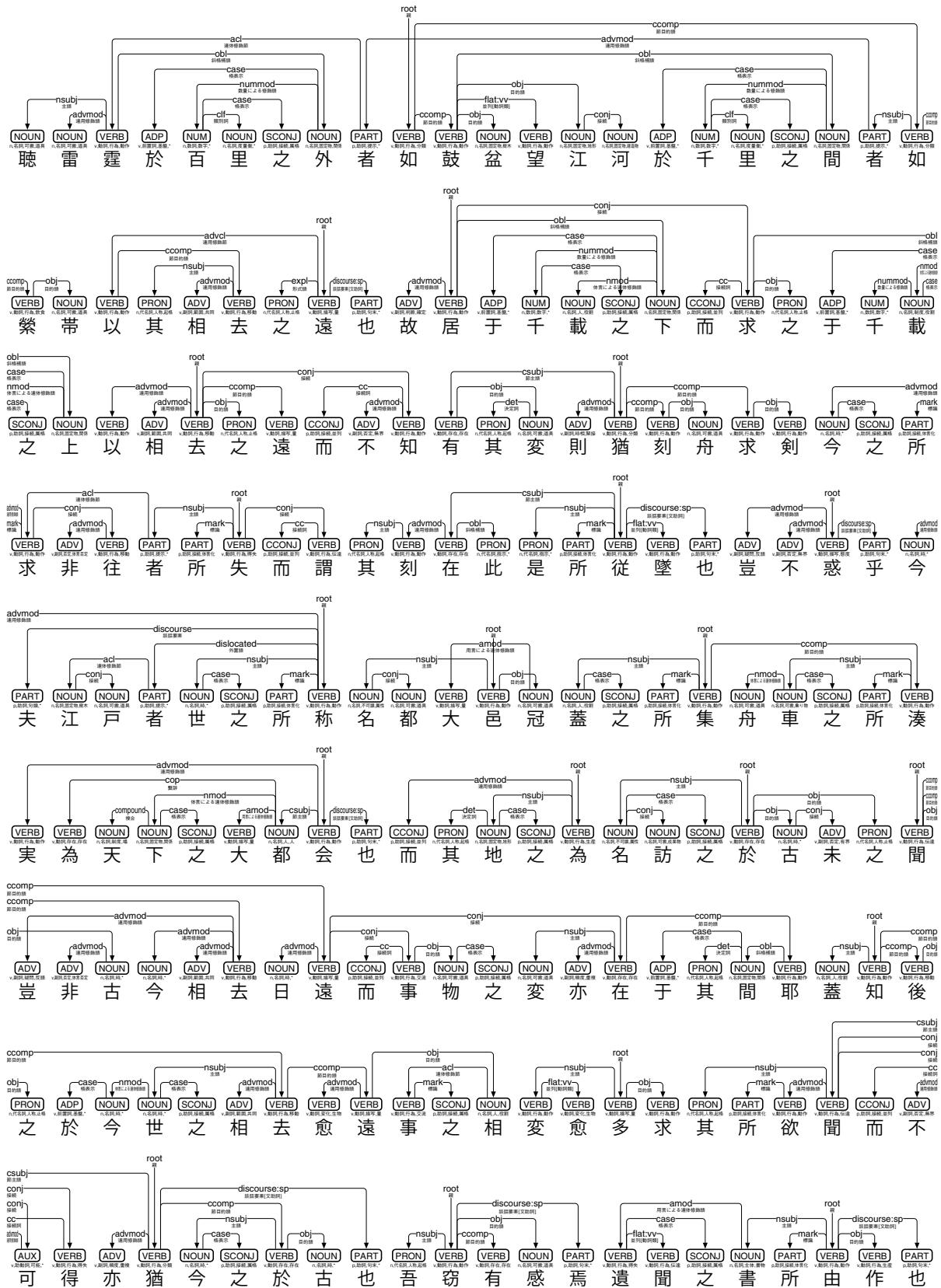


図 53: 手法Ⓐの処理結果(2017年)

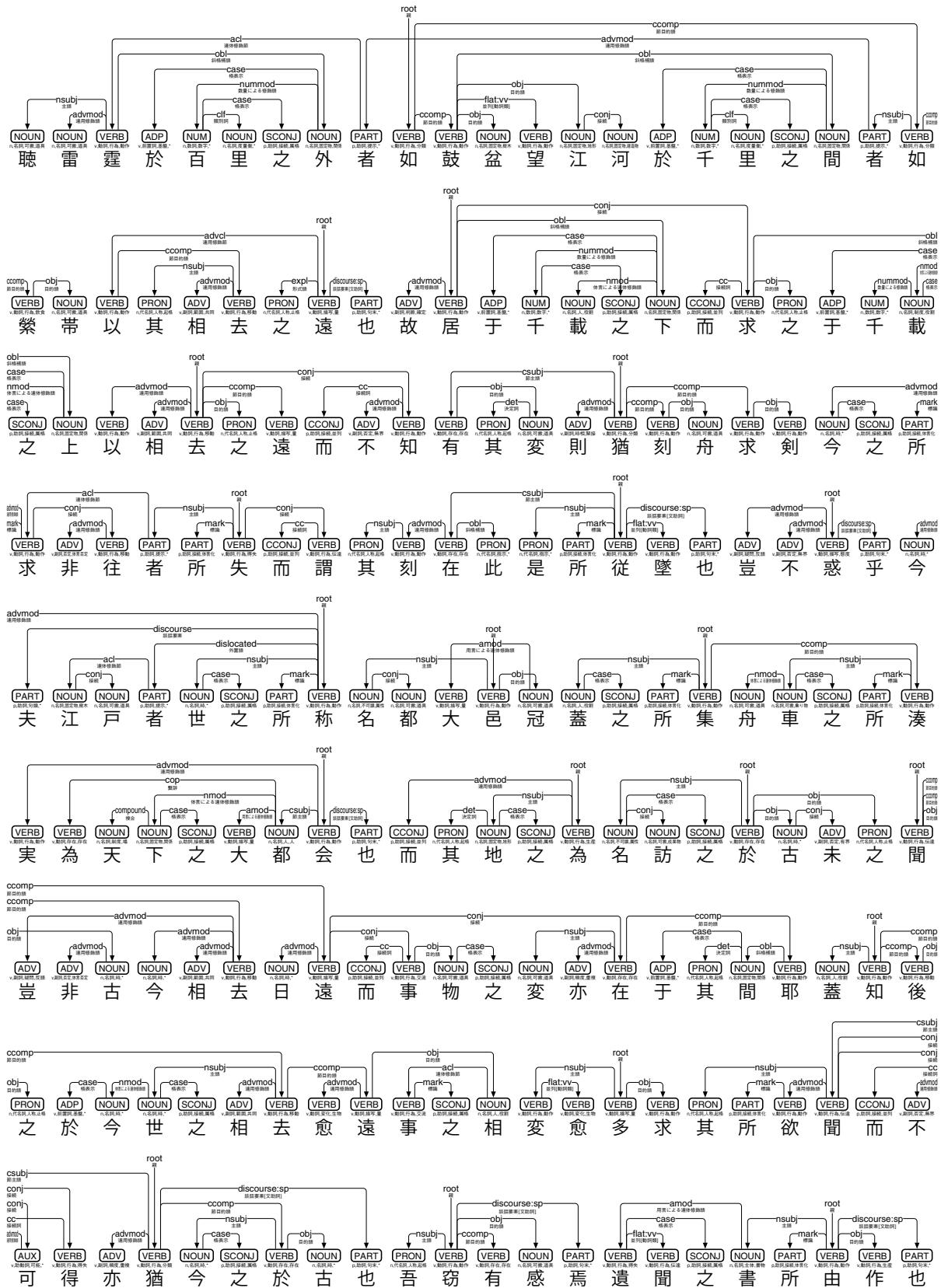


図 54: 手法②の処理結果(2017年)

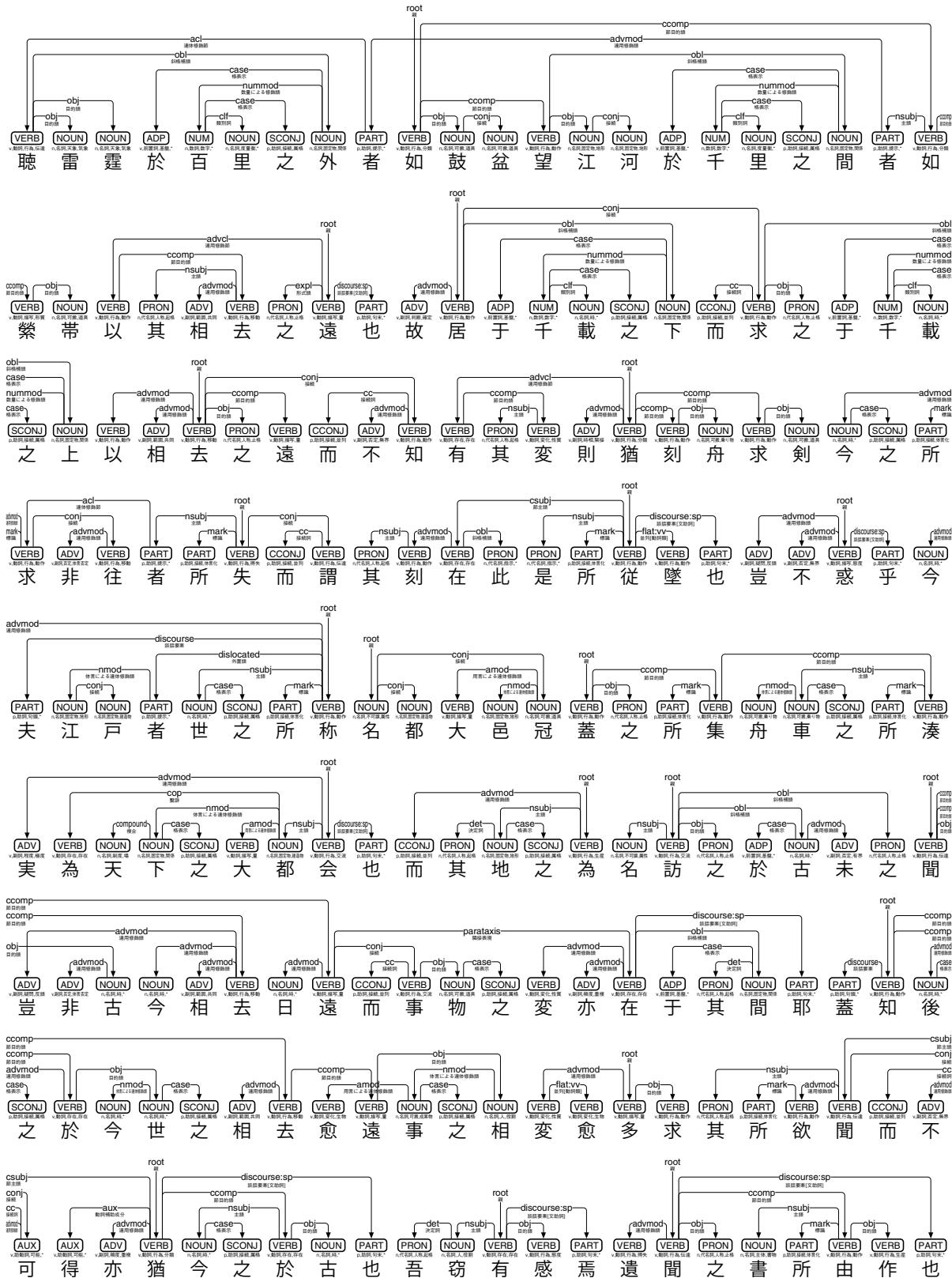


図 55: 手法④の処理結果(2017年)

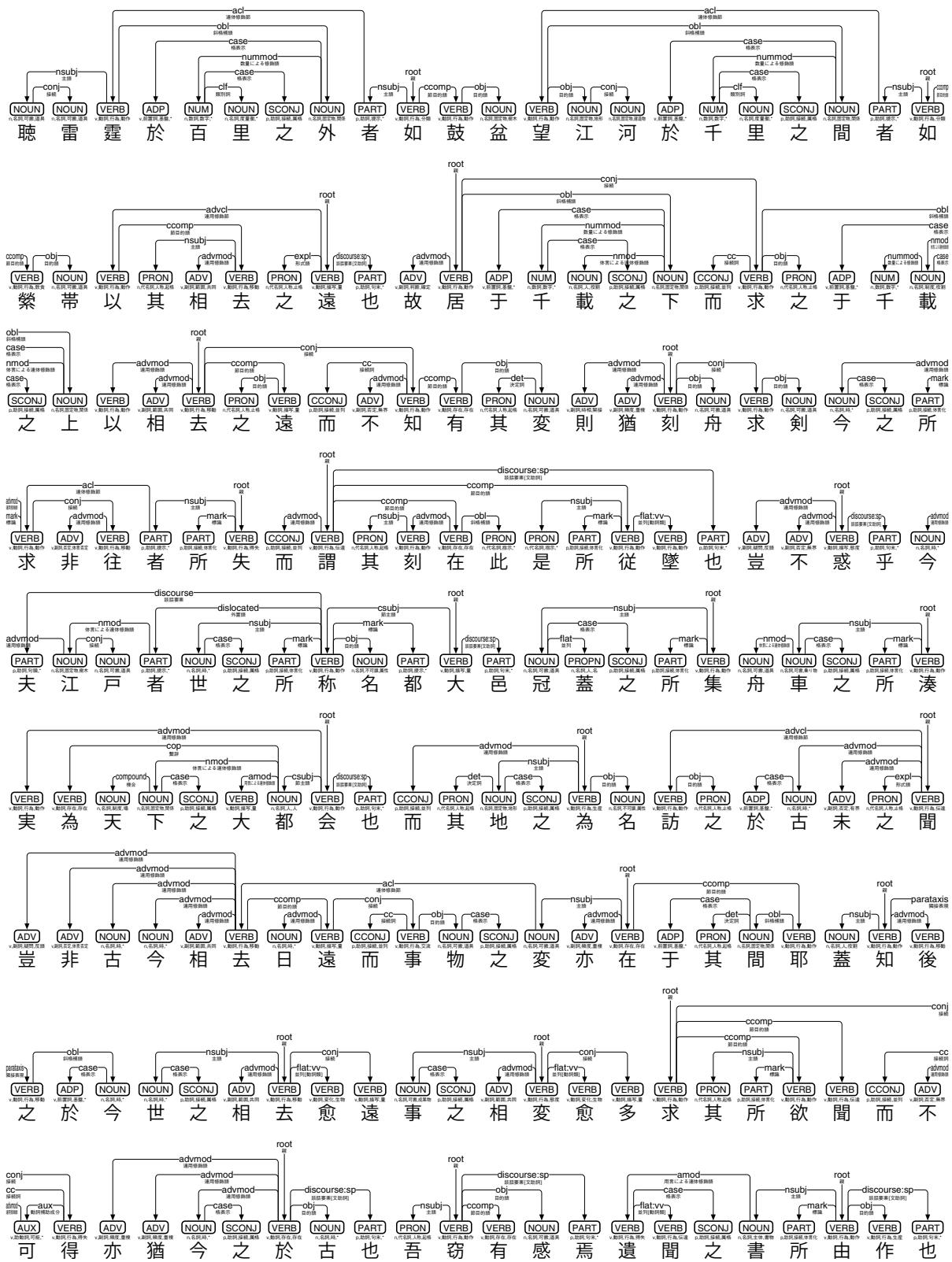


図 56: 手法Aの処理結果(2017年)

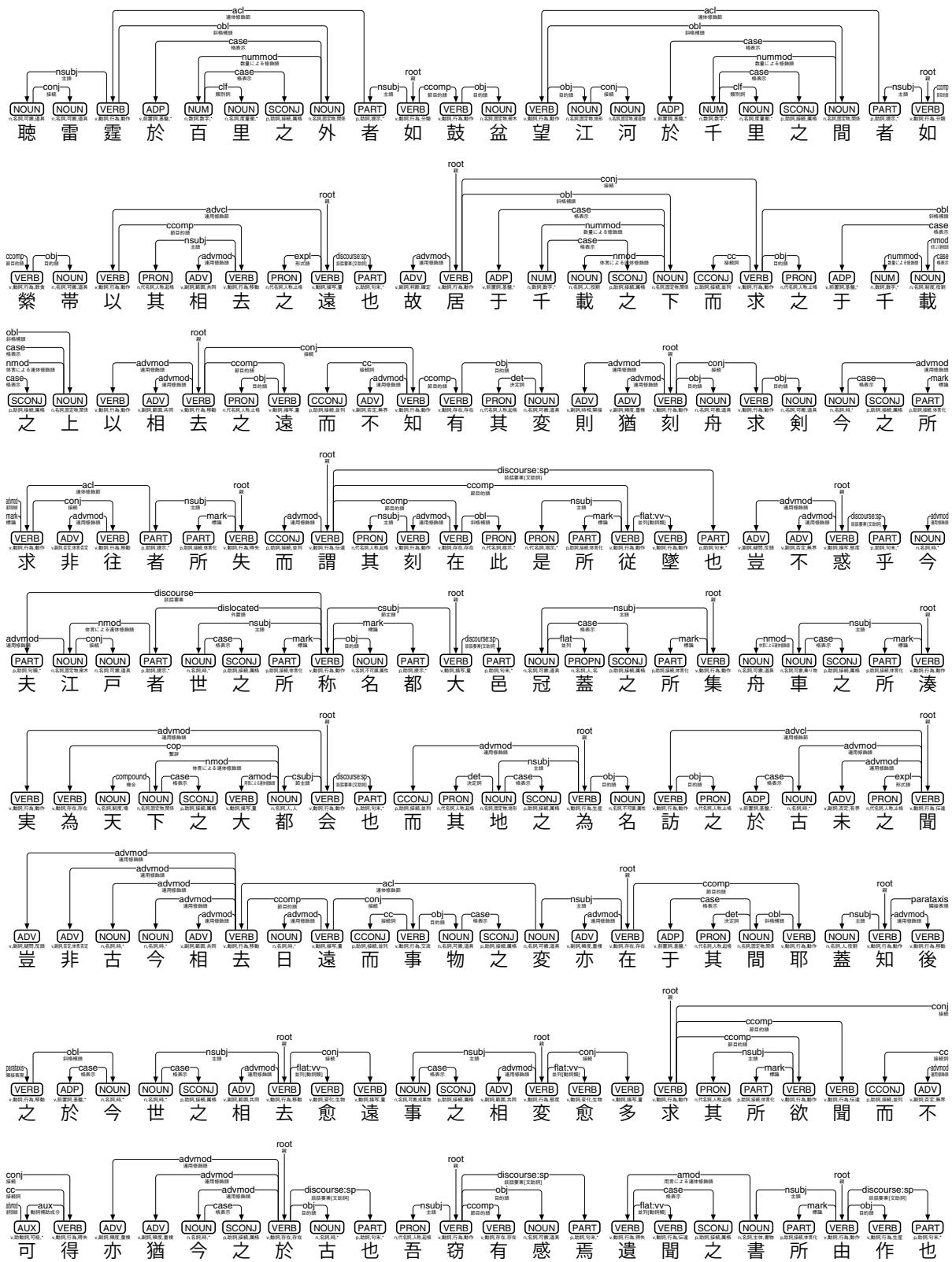


図 57: 手法Bの処理結果(2017年)

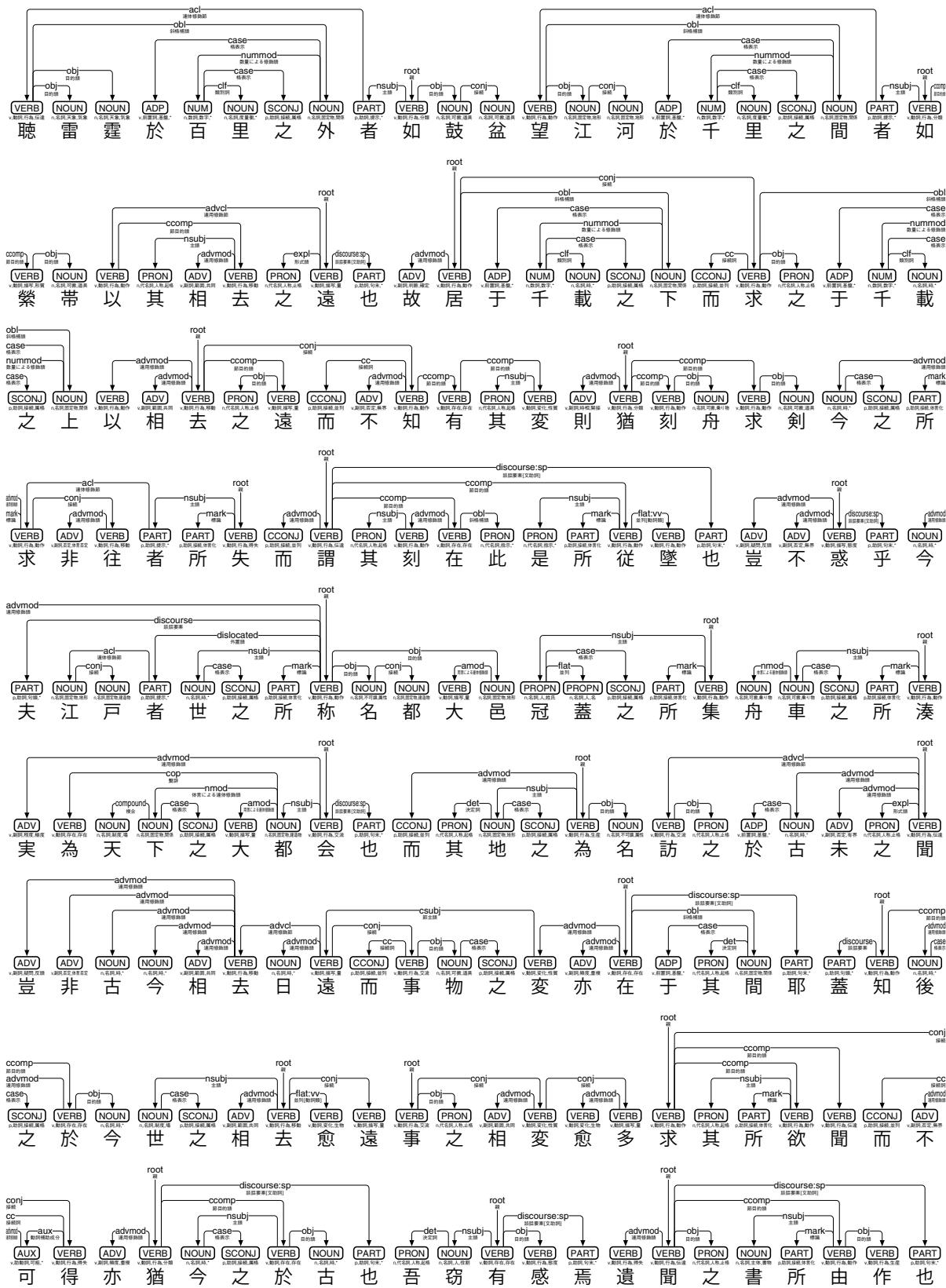


図 58: 手法Cの処理結果(2017年)

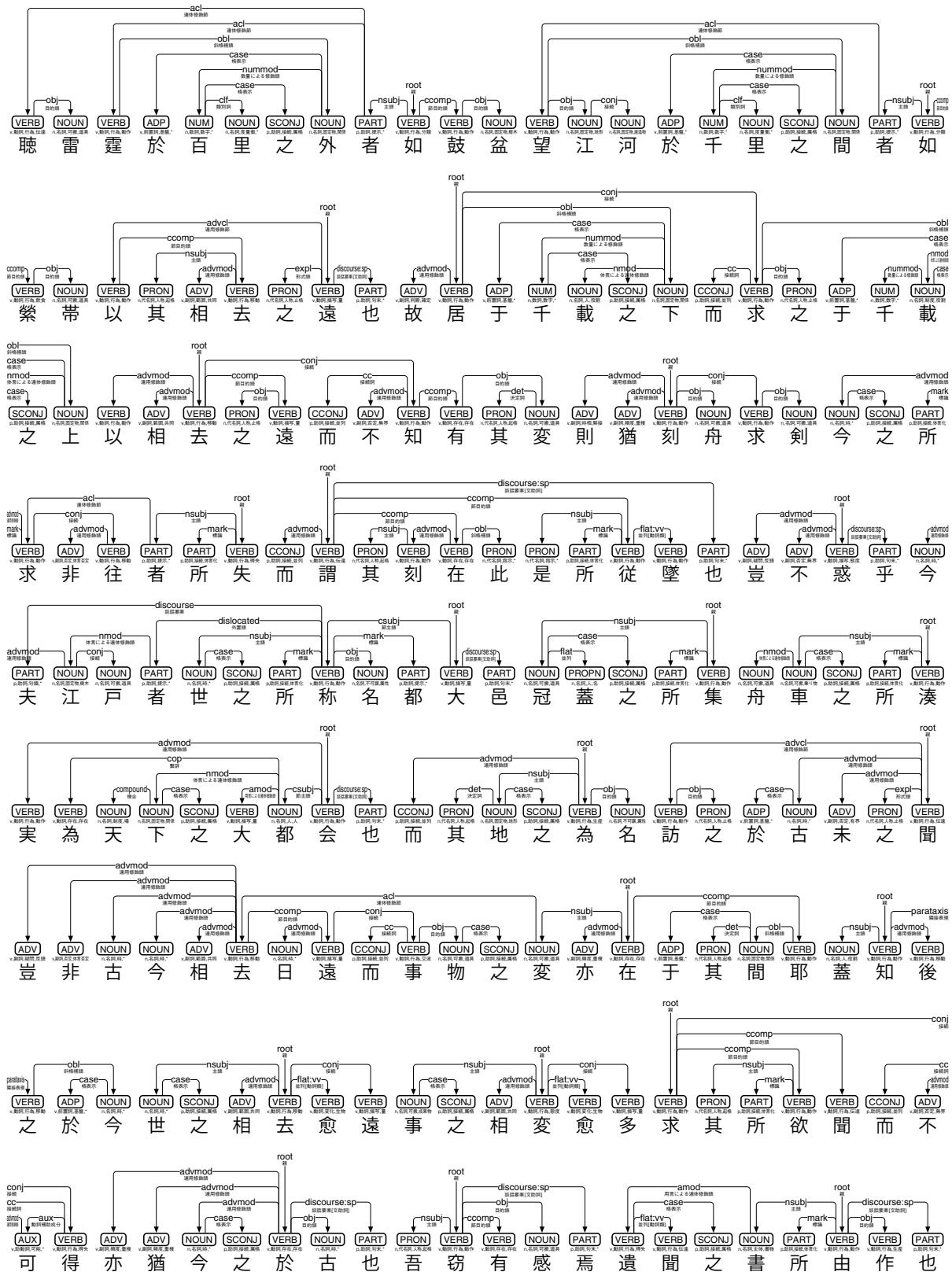


図 59: 手法Aの処理結果(2017年)

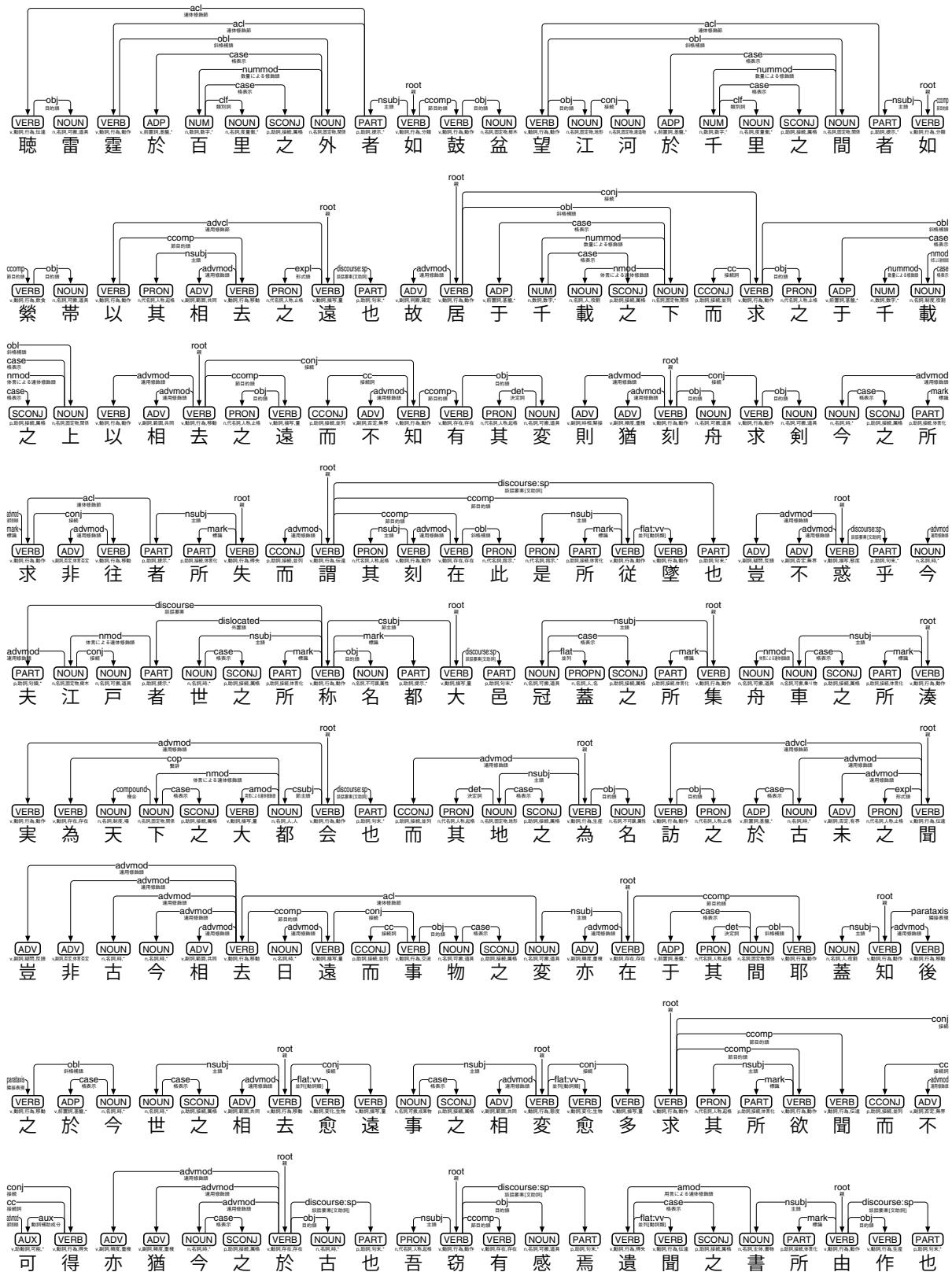


図 60: 手法Bの処理結果(2017年)

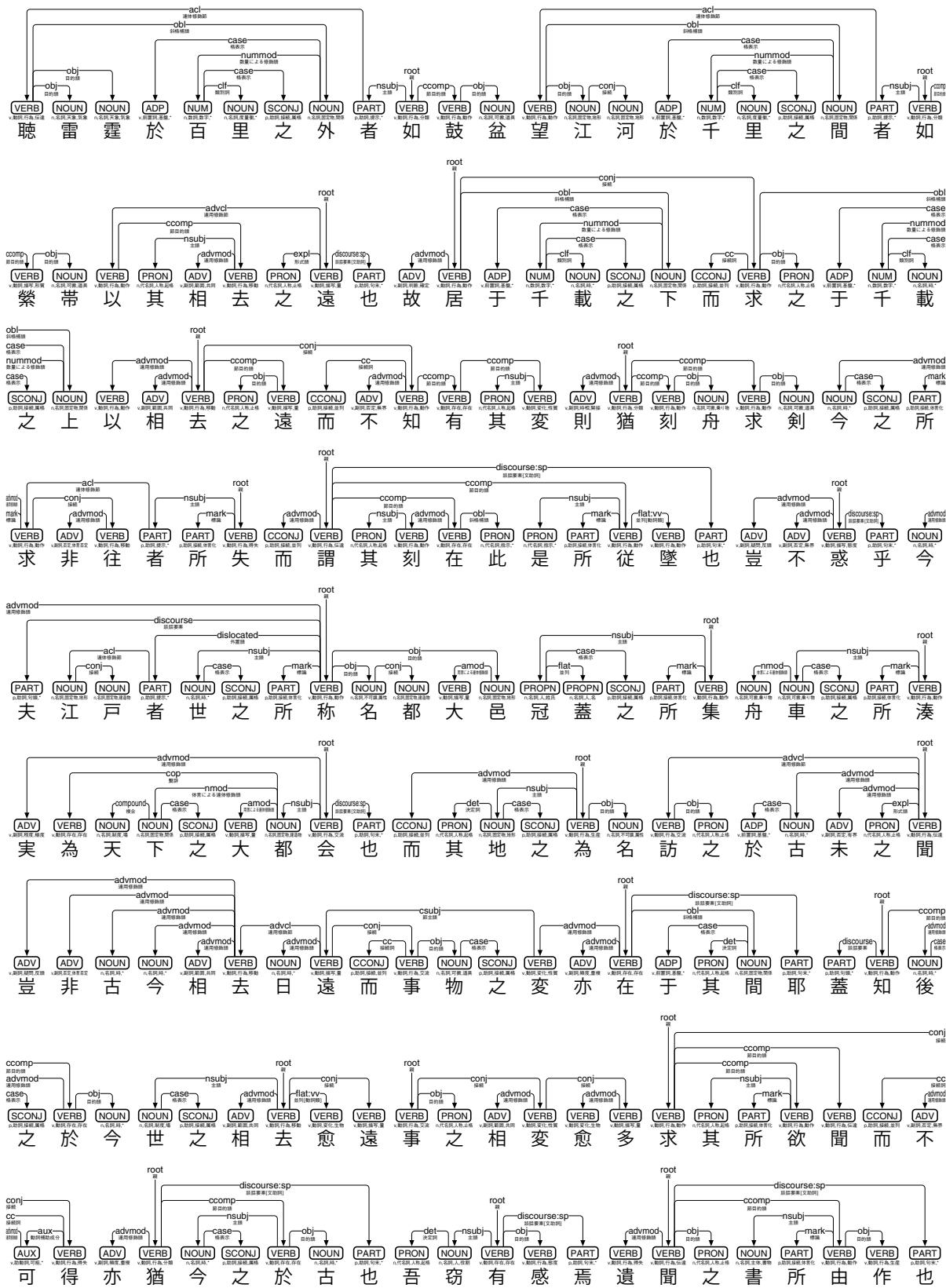


図 61: 手法Cの処理結果(2017年)

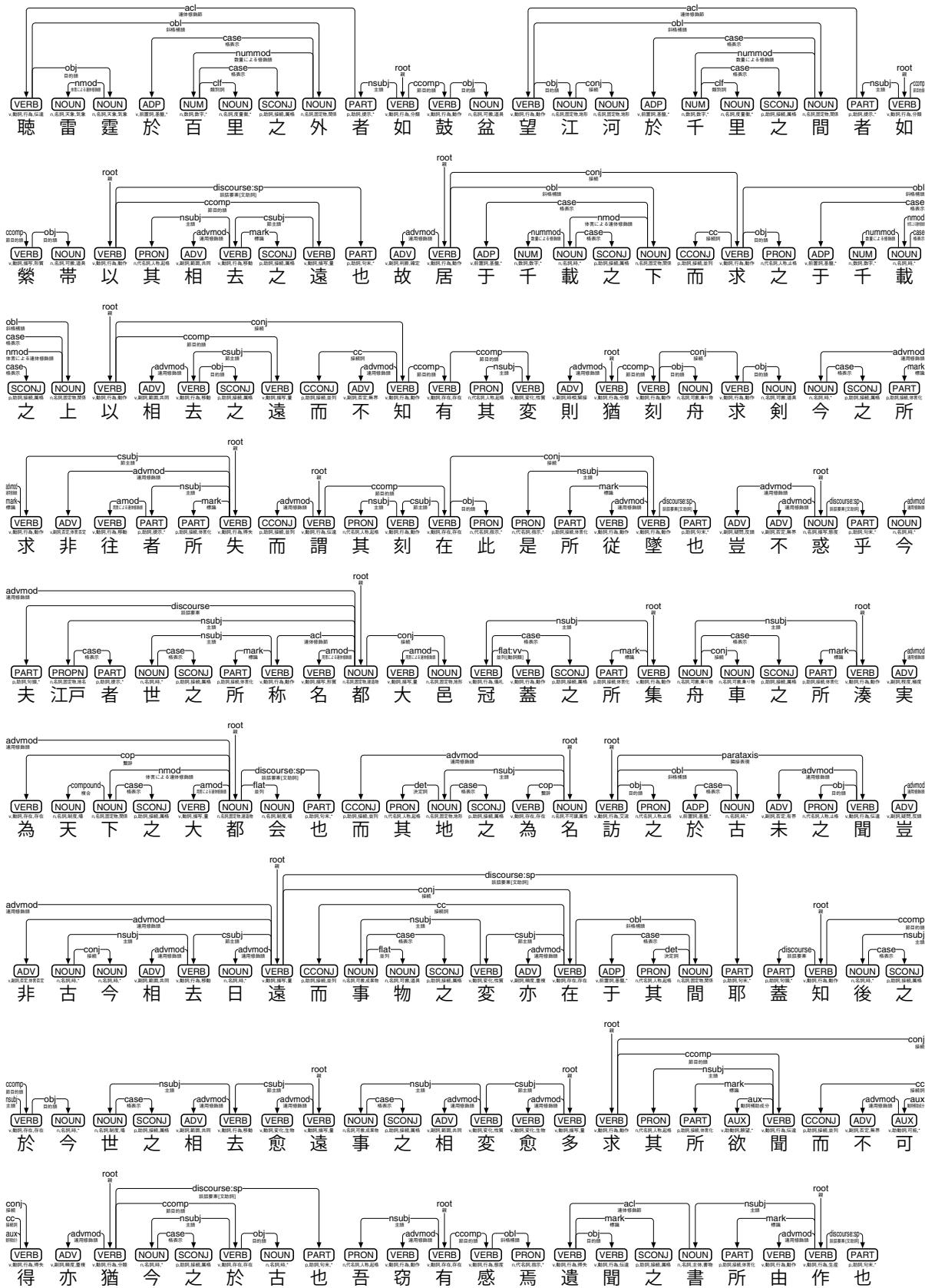


図 62: 手作業で作成した「正解」UD (2017 年)

者乎。ト

A 荷宇ハシマレテ生ニシテ十月而喪ウシナフ其母ヲ及ビ(1)レ有ルニ知シ即チ時カナシム時モ念メヒテ母不レ置カイヨイ弥久シクシテ

B 弥篤アツシ哀ハシマレ其身不能ハシマレ一日事乎母也ヲ哀母之言語動作亦未ダル

C 能ハシマレ識シル也。

D 繼ギテ此而得ル見ル也ト於レ是ニ即シテ夢ミル所レ見ル為ツカル之ガ圖ヲ此圖ハ吾不ル之ヲ見ル也。

E 使ムル我ヲシテ至リテ今ニ日ニ乃チ得ル見ル也ト母ヨ又タ何ソ去ルコト我ヲ之ヤカル速ナル也。母ヨ其レ可ケン使ム我ヲシテ母ヨ胡ナン為スレゾ乎。

F 知ル其ノ為ル母也既ニ覺メ(イ)乃チ噭けい然トシテ以テ哭ハシマレ曰ハク此真吾母ガ也。母ヨ母ヨ胡ナン為スレゾ乎。

今之図吾見ルニ之ヲ則チ其ノ夢ミル母ヲ之ナル境ナ而已。

余因語リテ之ニ曰ハク夫レ人精誠所レ感ズル無キハ幽注5明死生之隔テ此理レ之

可ク信ス不レ誣シヒ者ナリ況ンヤ子之於ケル親ニ其ノ喘ゼン息そき呼吸モ相通ジ本ヨリ無キヲ有ル間ヘダツル之ヲ

図 63: 大学入試センター試験『国語』(2016年1月16日)第4問本文

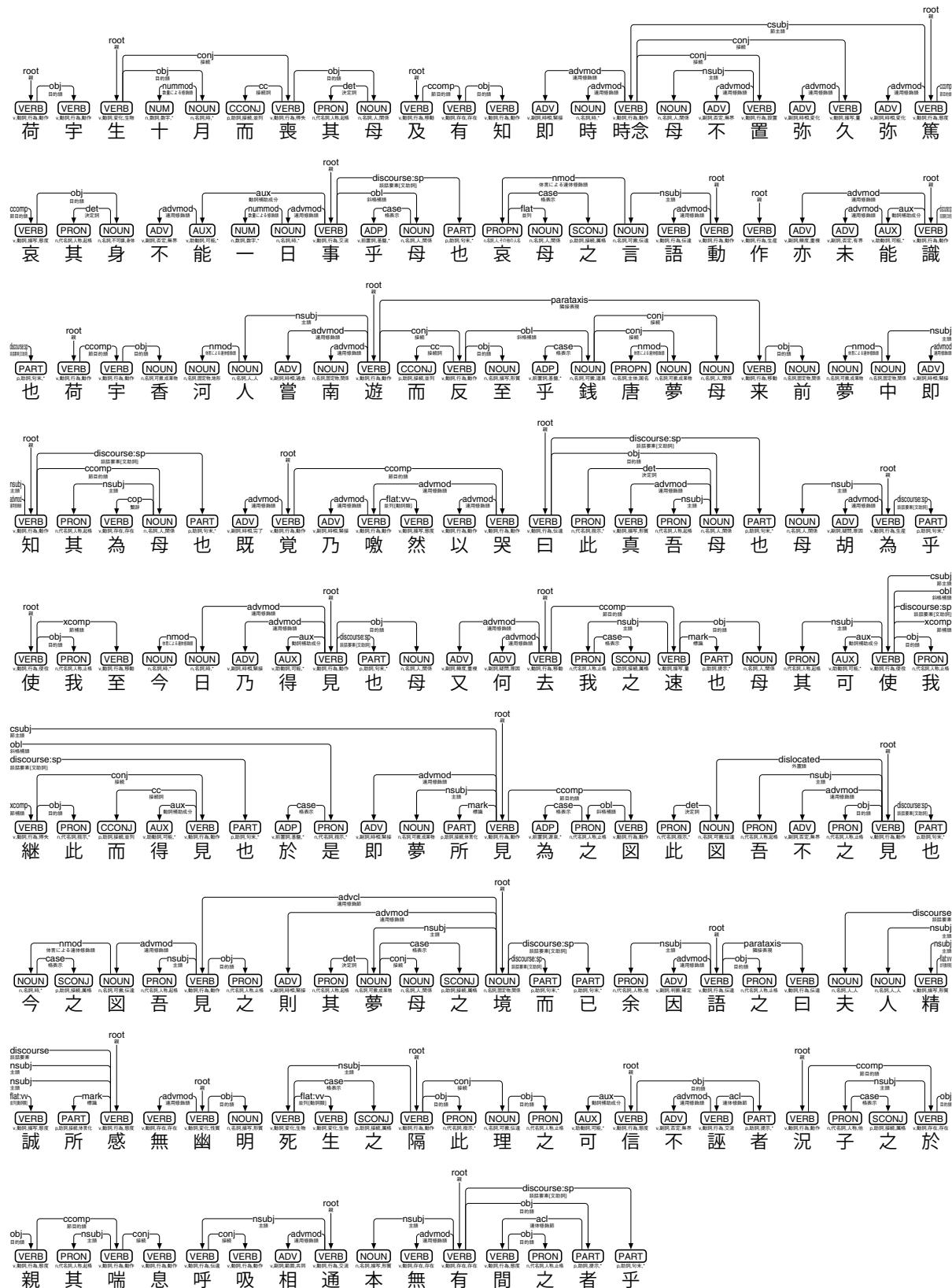


図 64: 手法①の処理結果(2016 年)

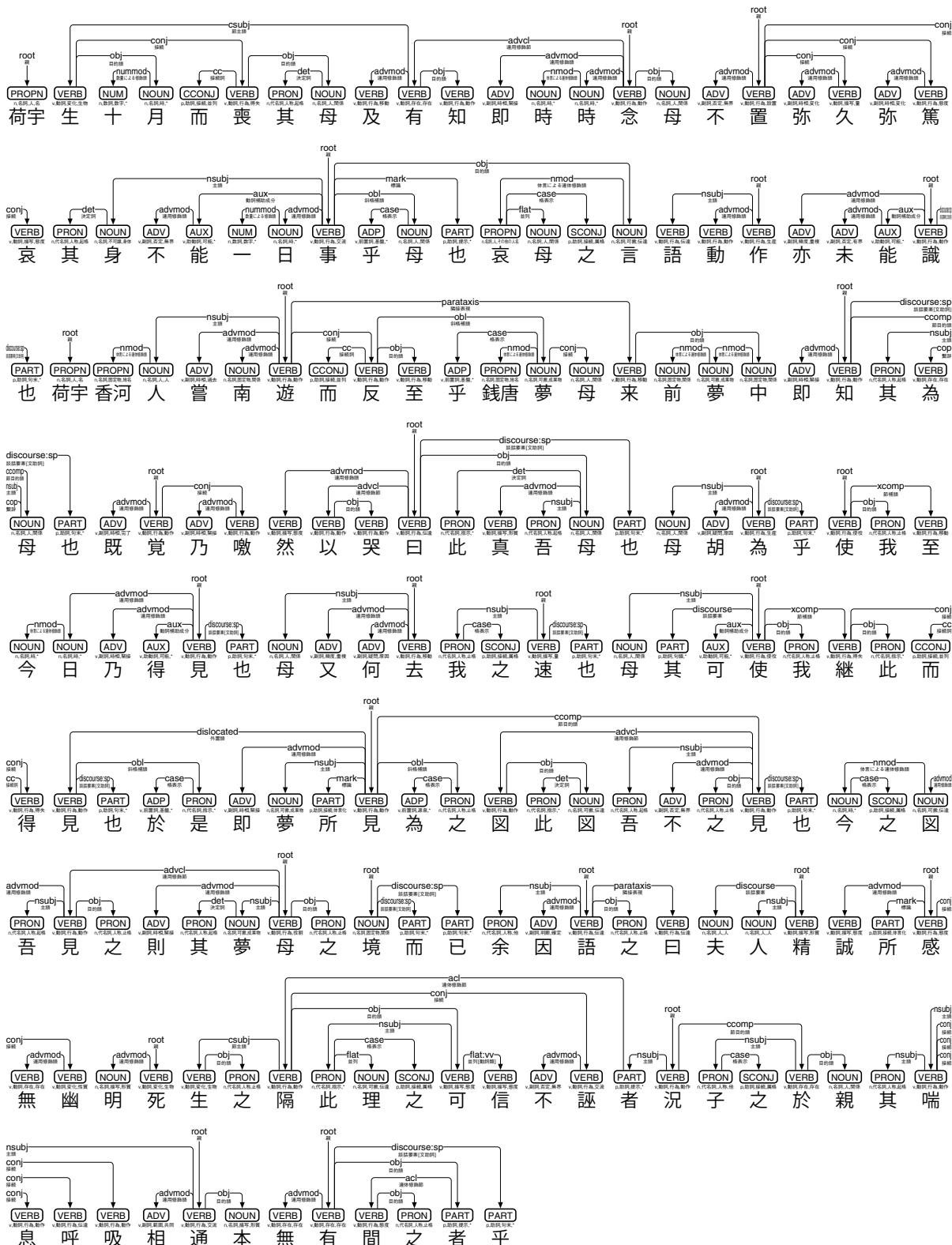


図 65: 手法②の処理結果(2016年)

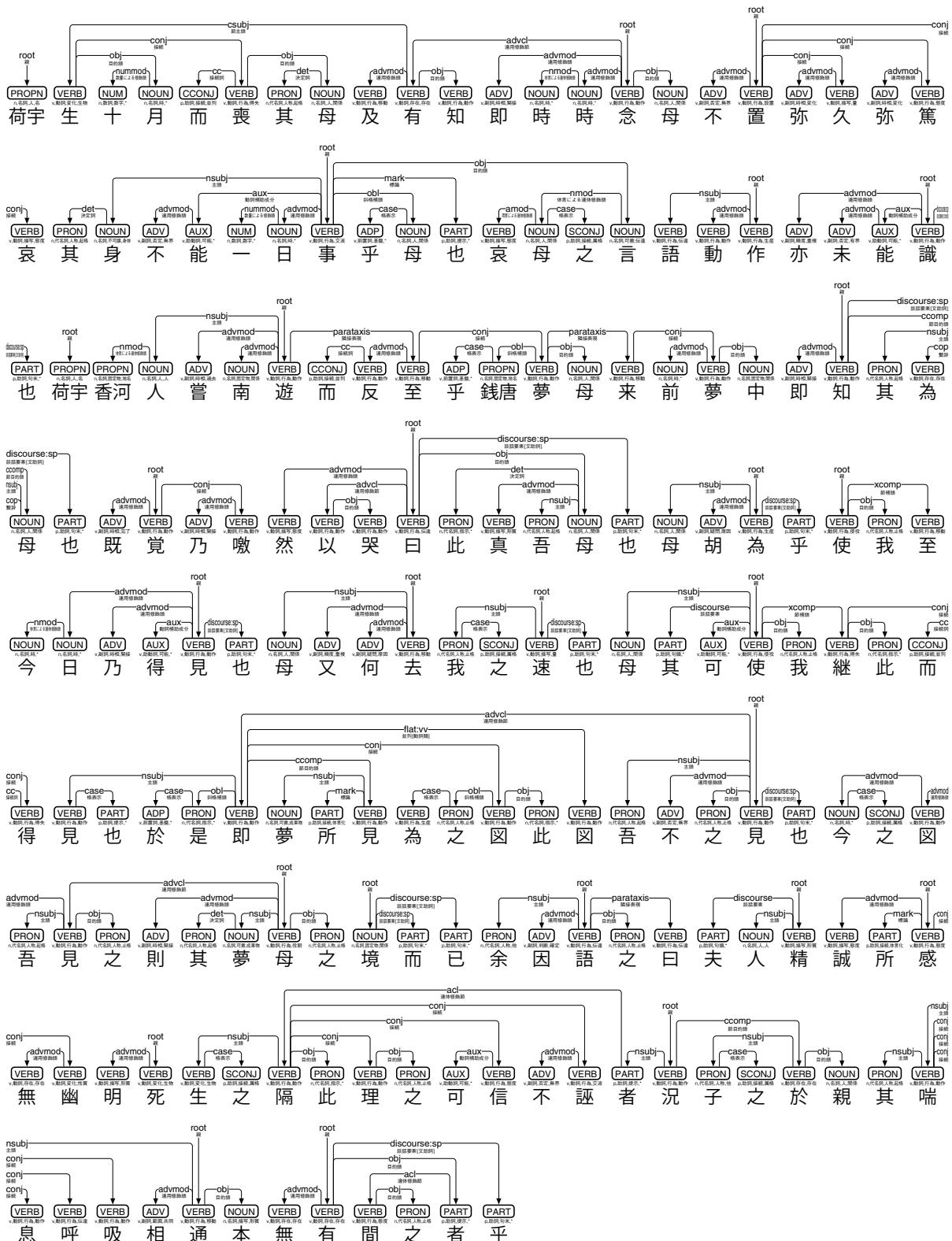


図 66: 手法③の処理結果(2016年)

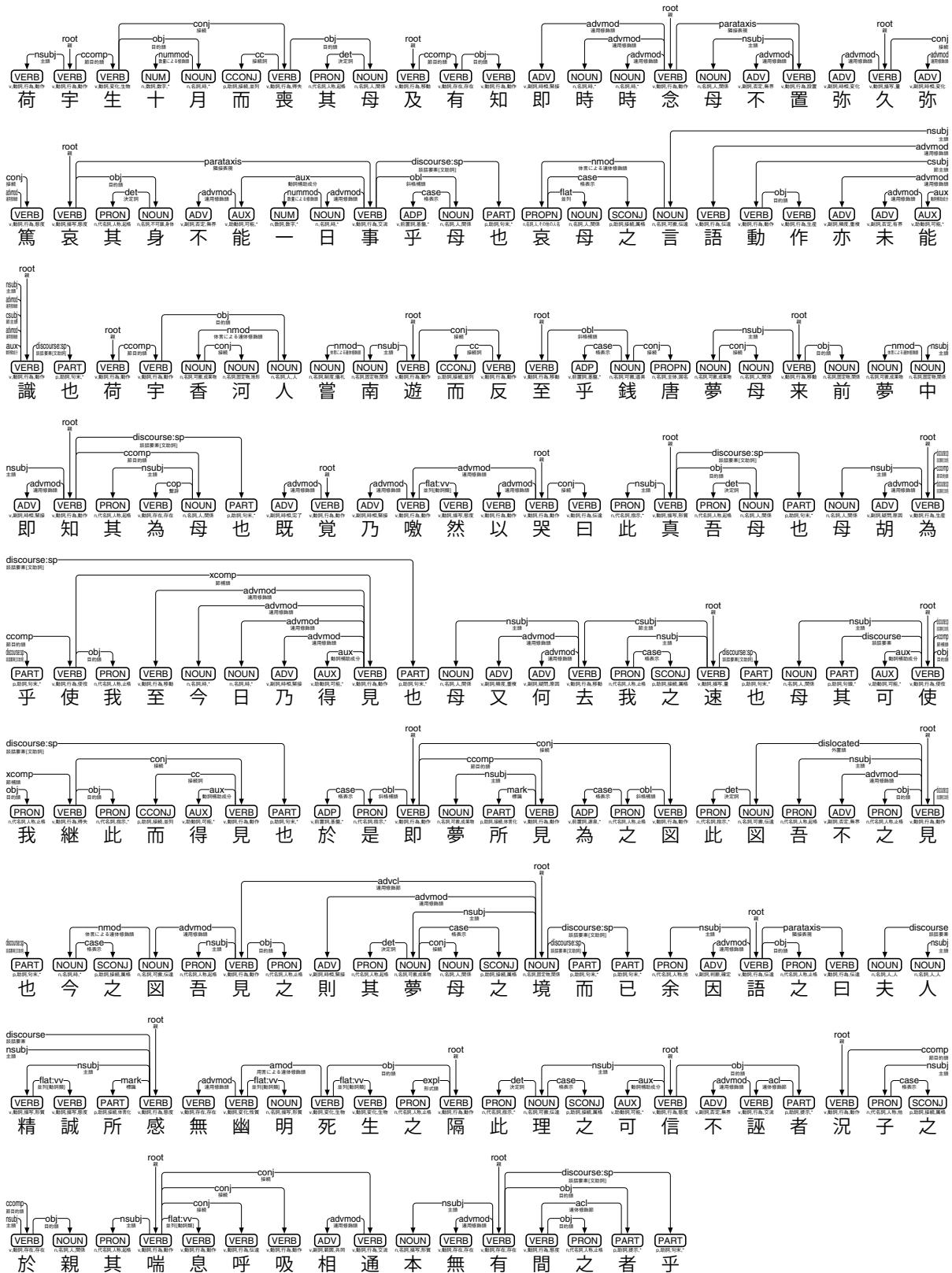


図 67: 手法①の処理結果(2016年)

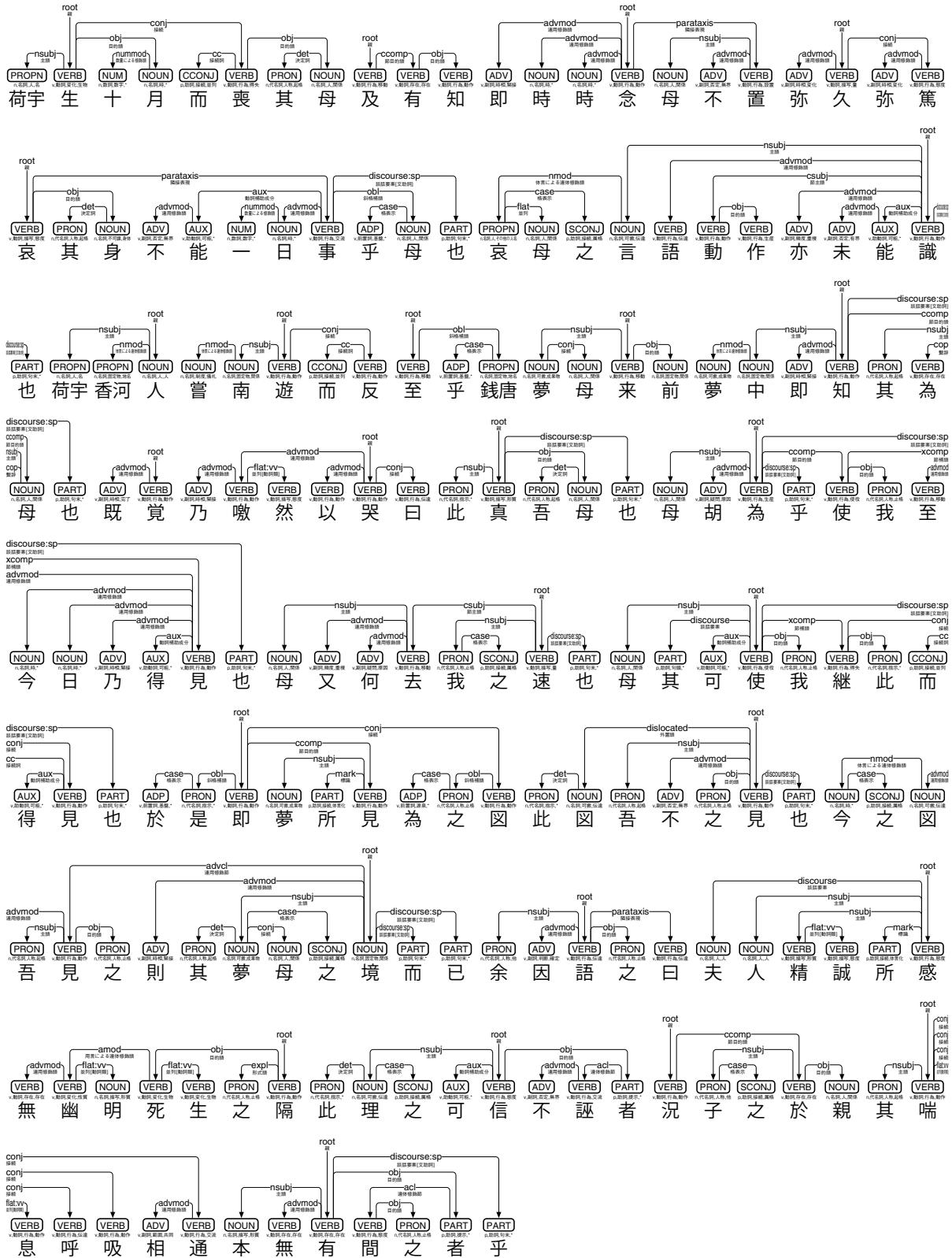


図 68: 手法[2]の処理結果(2016 年)

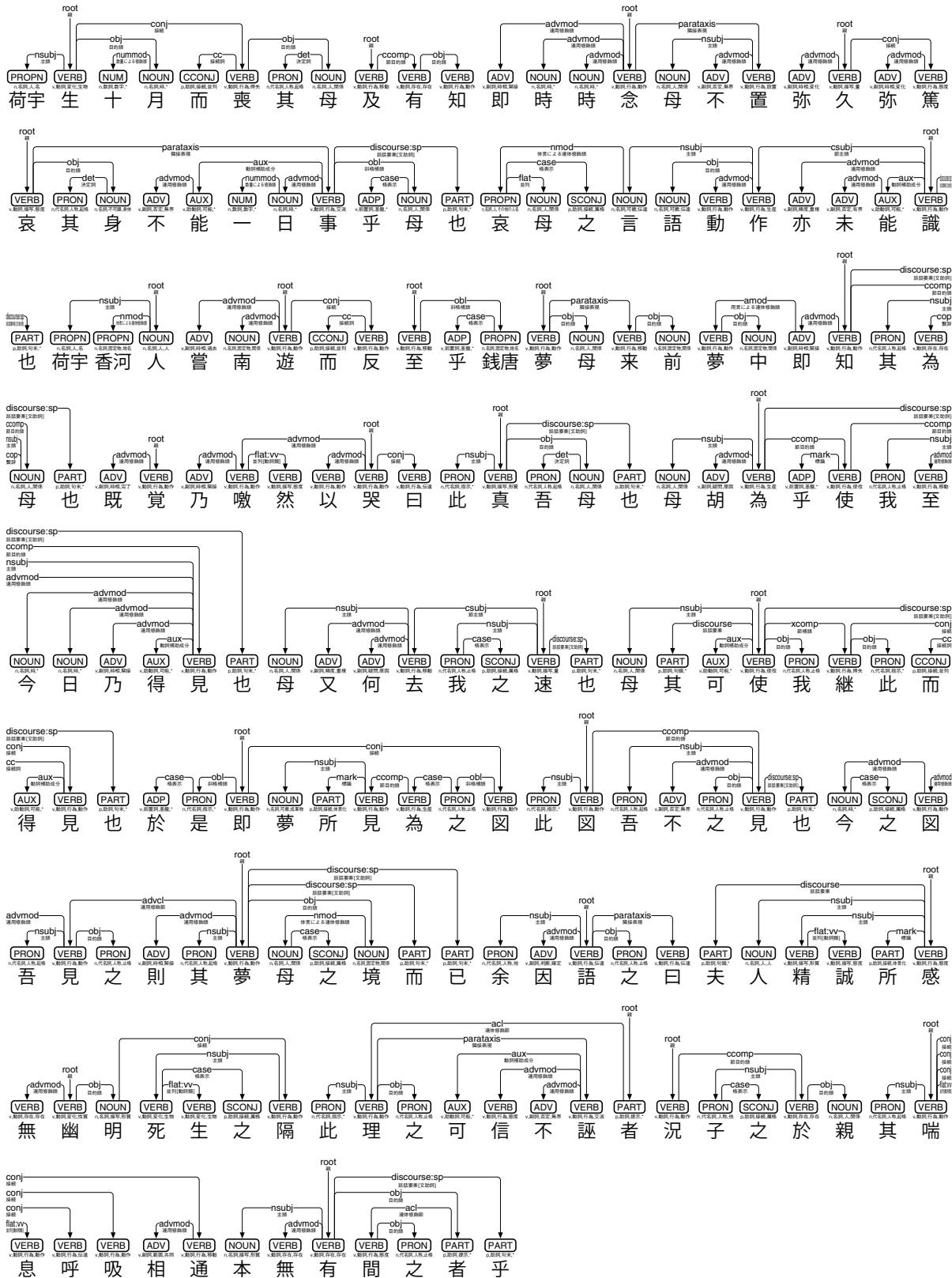


図 69: 手法③の処理結果(2016年)

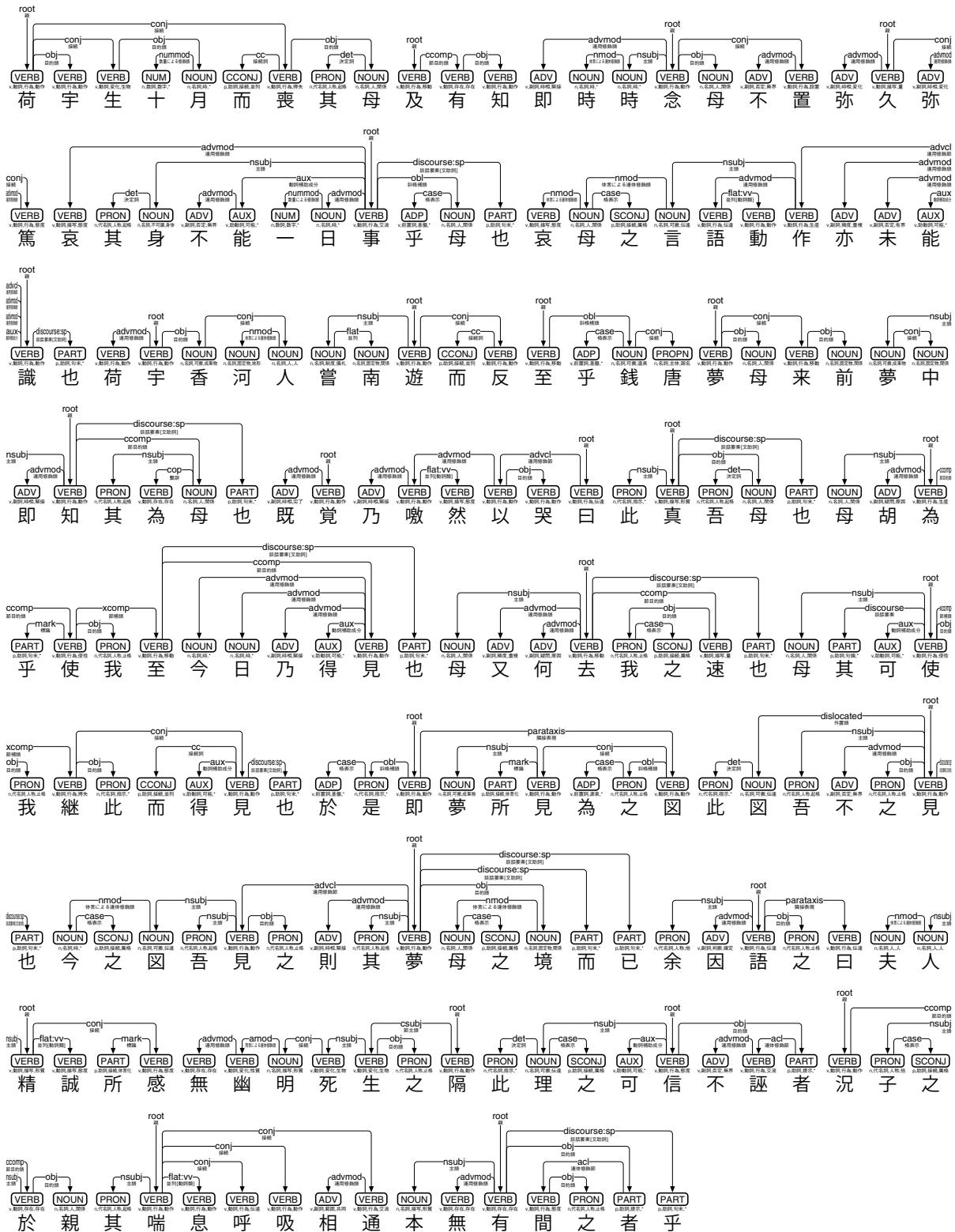


図 70: 手法①の処理結果(2016年)

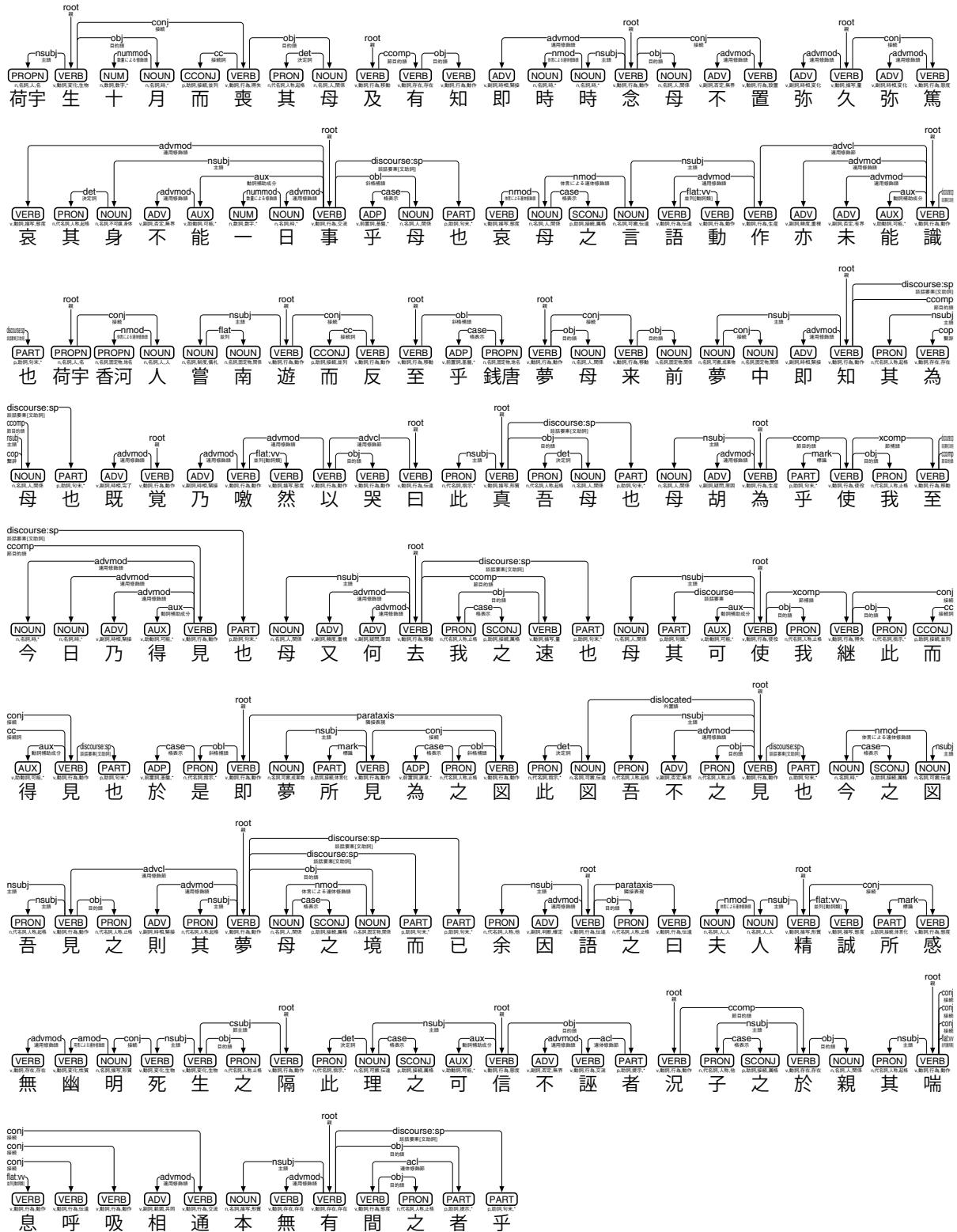


図 71: 手法②の処理結果(2016年)

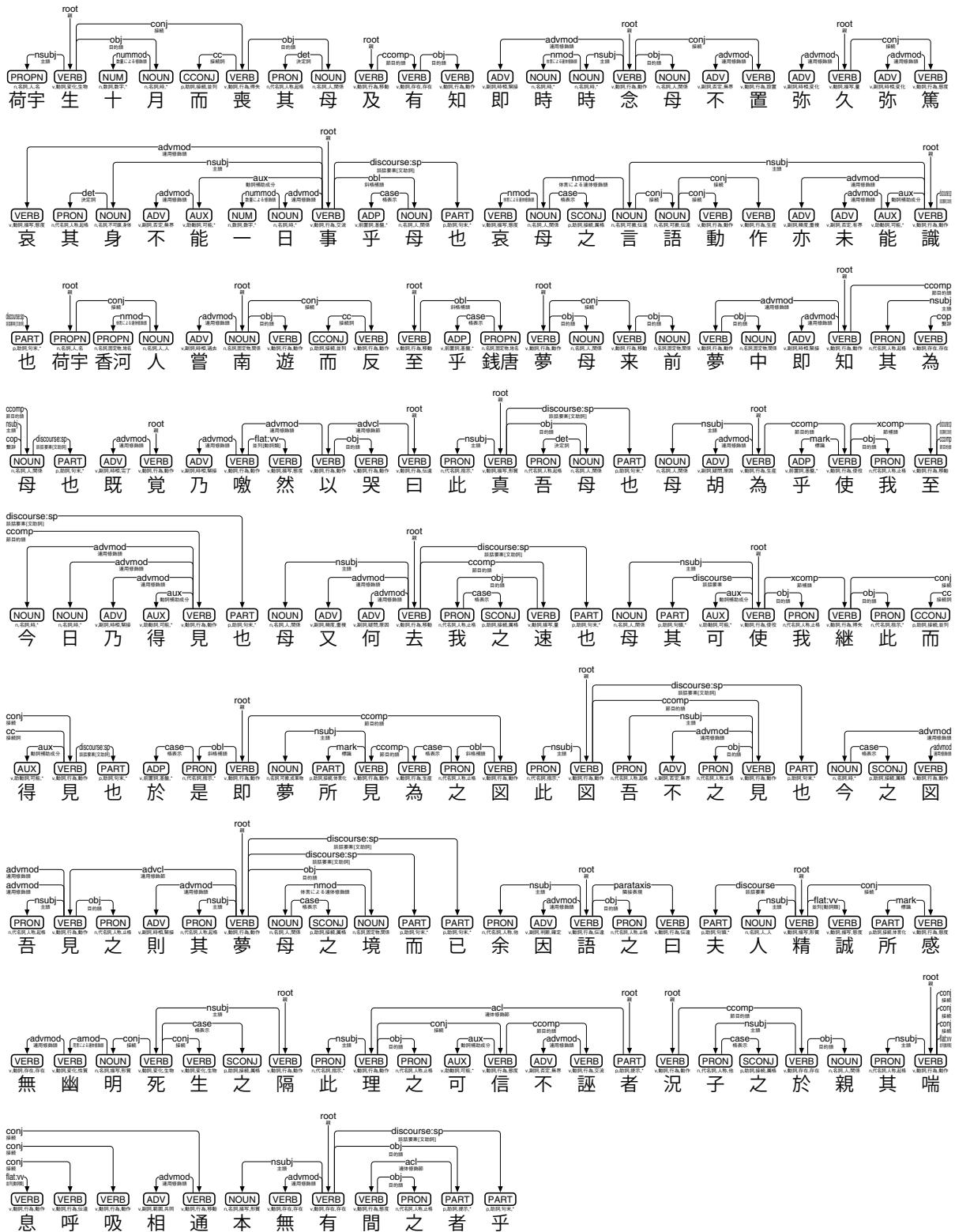


図 72: 手法③の処理結果(2016年)

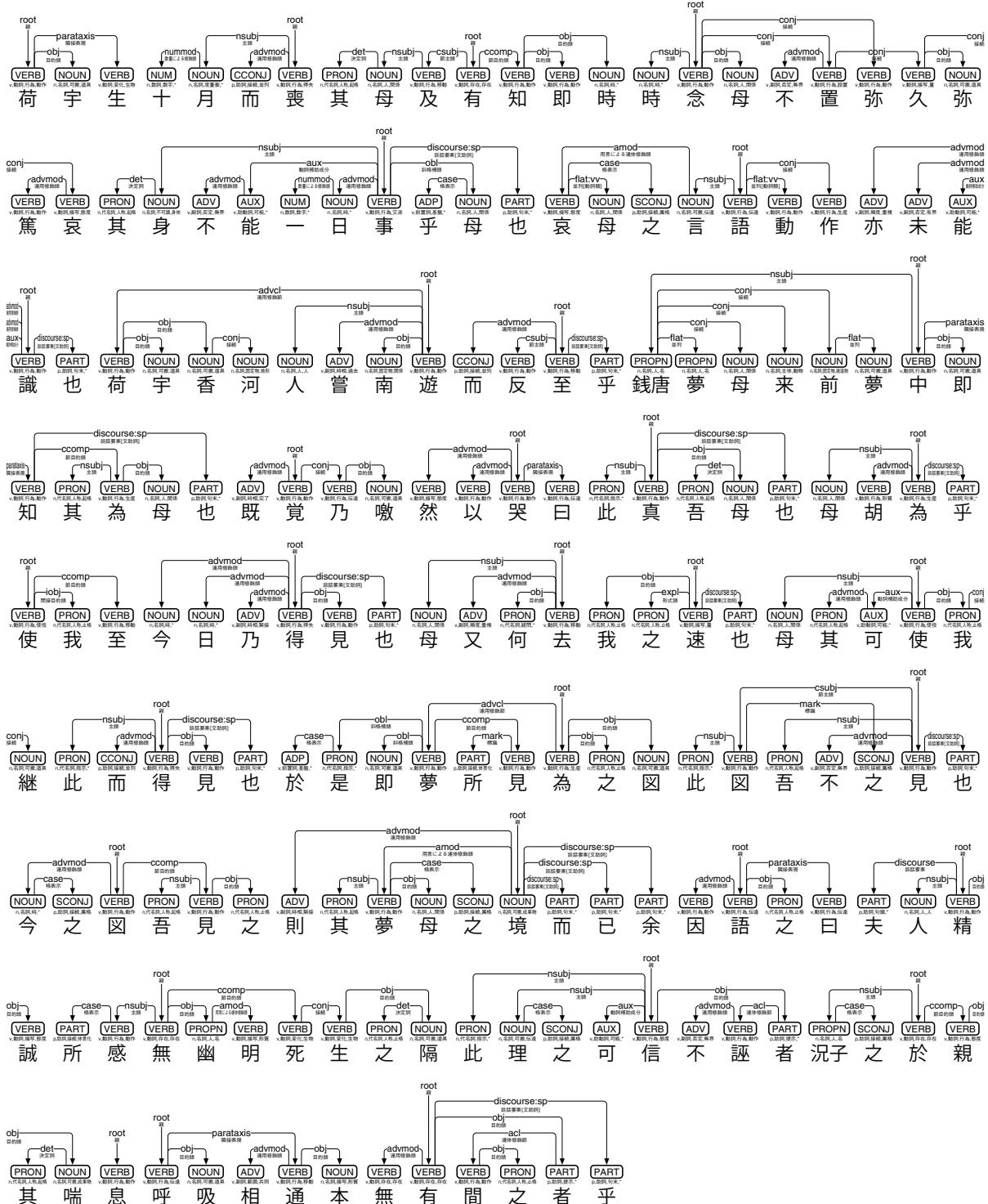


図 73: 手法Ⓐの処理結果(2016年)

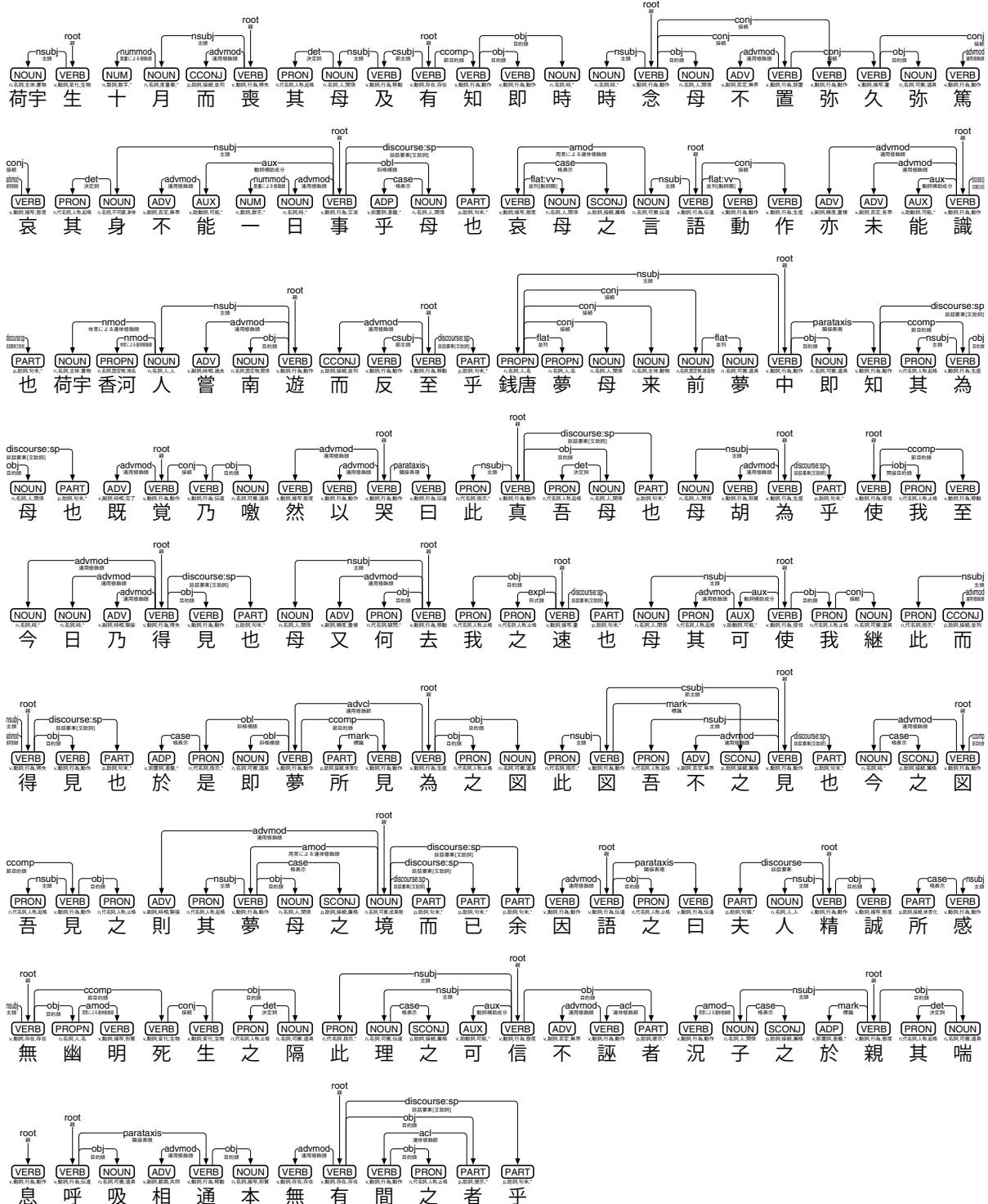


図 74: 手法②の処理結果(2016年)

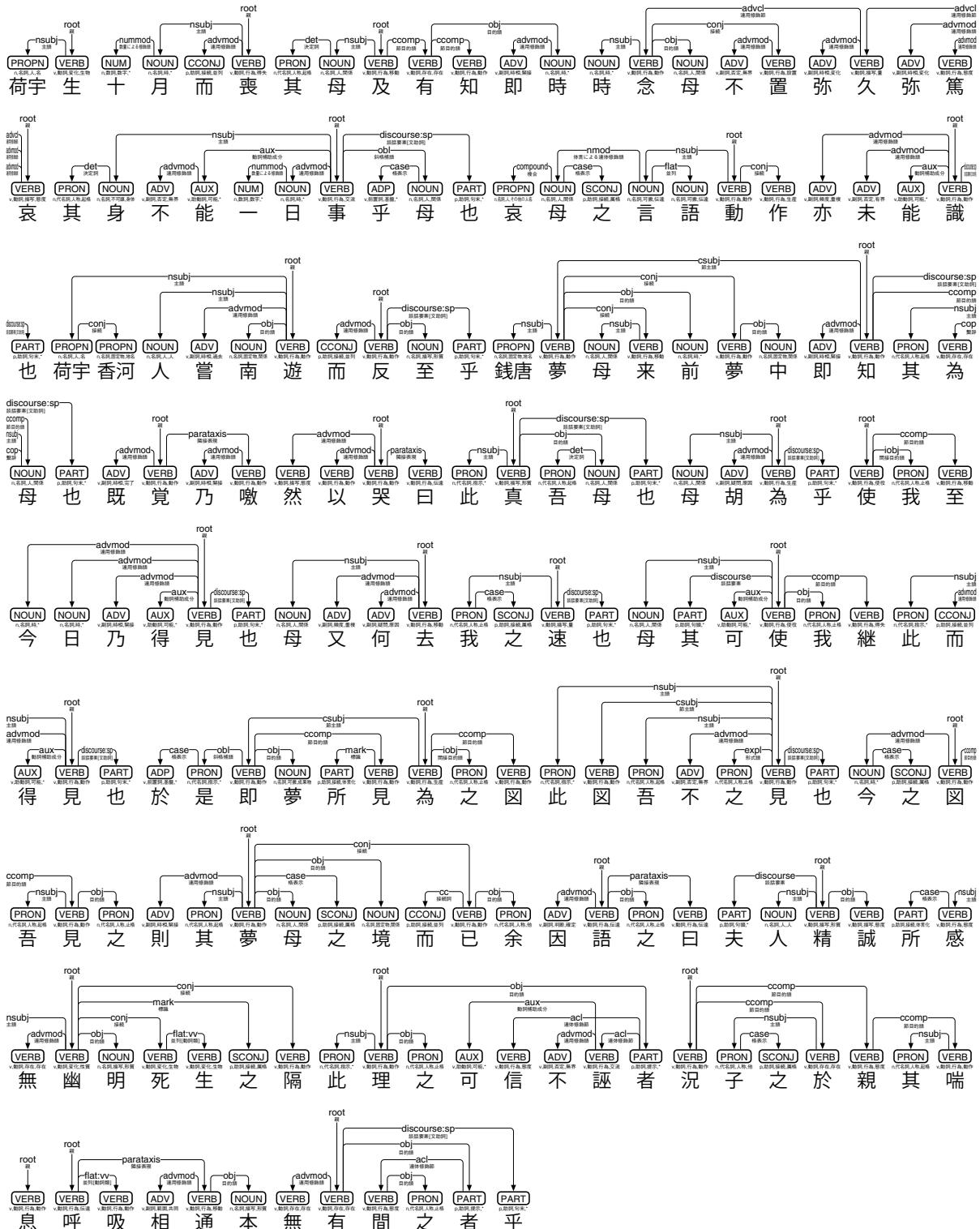


図 75: 手法③の処理結果(2016年)

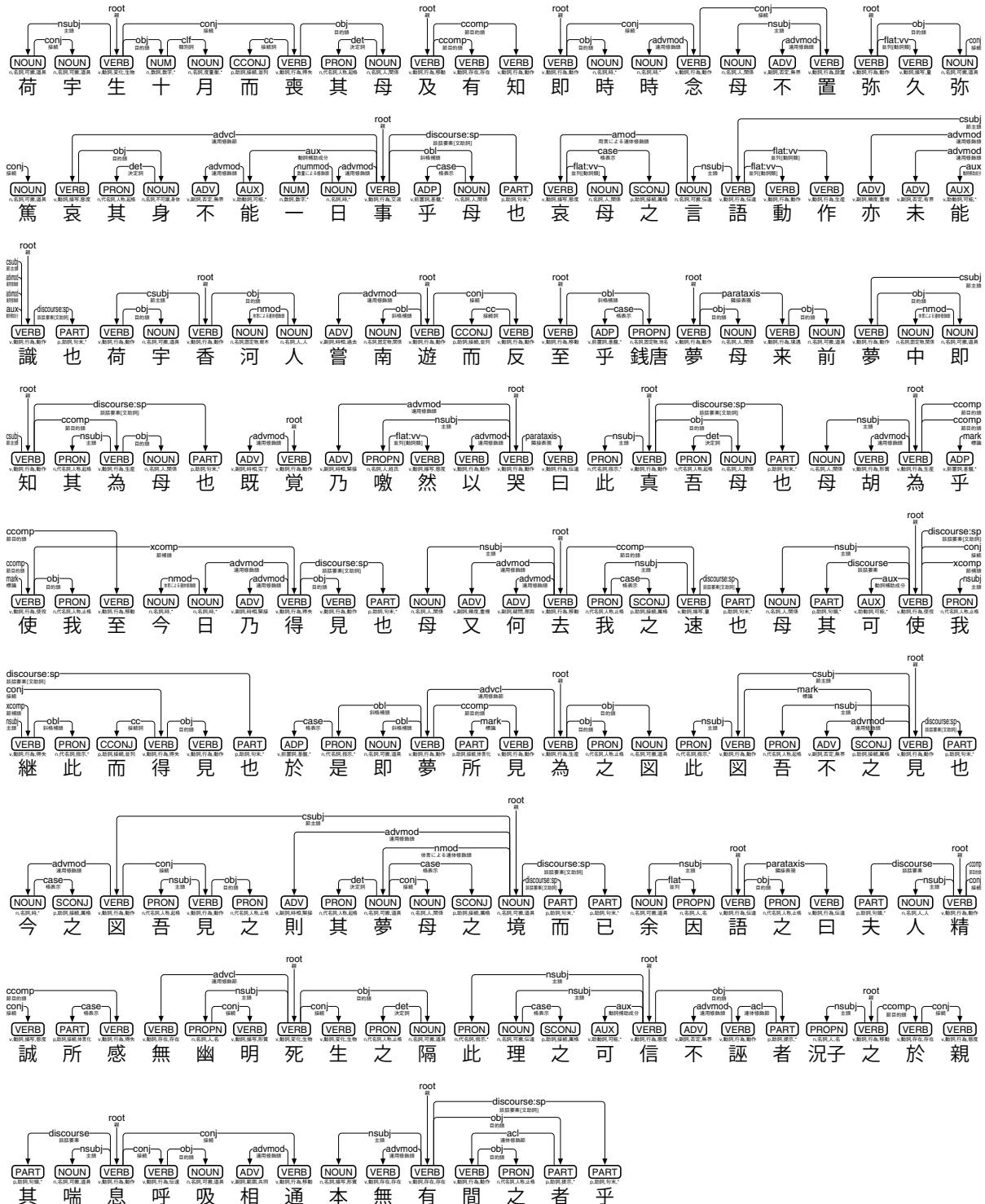


図 76: 手法Aの処理結果(2016年)

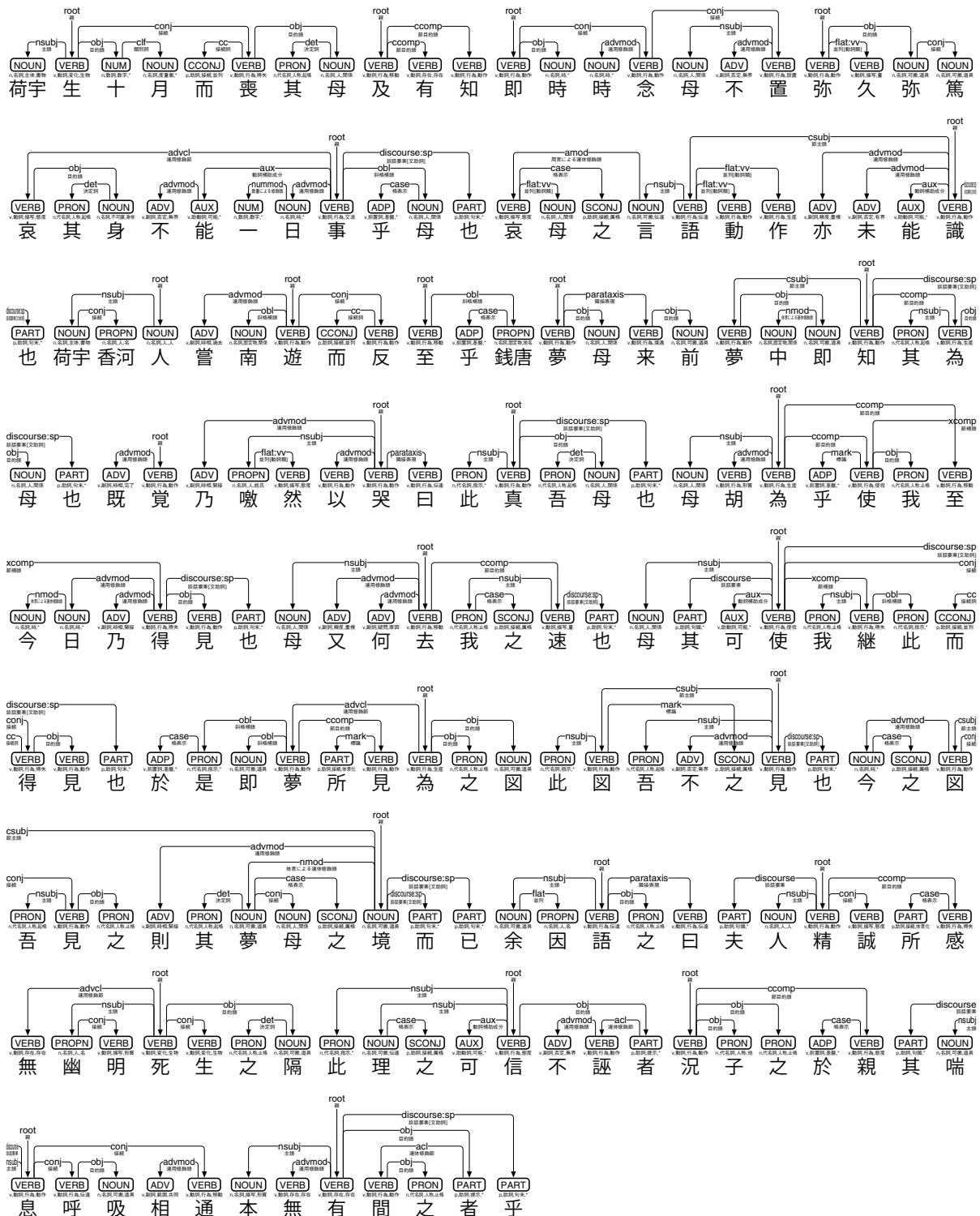


図 77: 手法Bの処理結果(2016年)

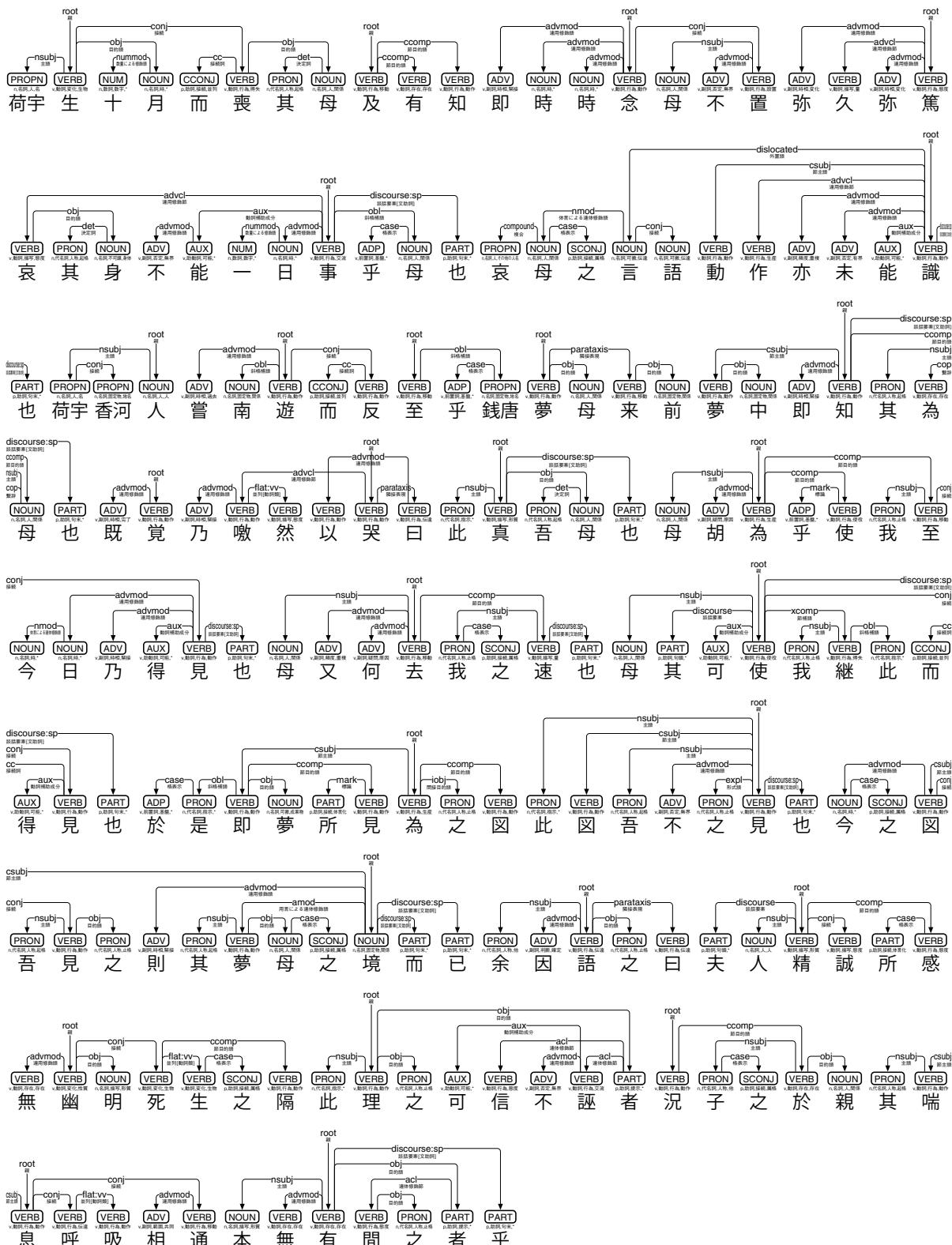


図 78: 手法Cの処理結果(2016年)

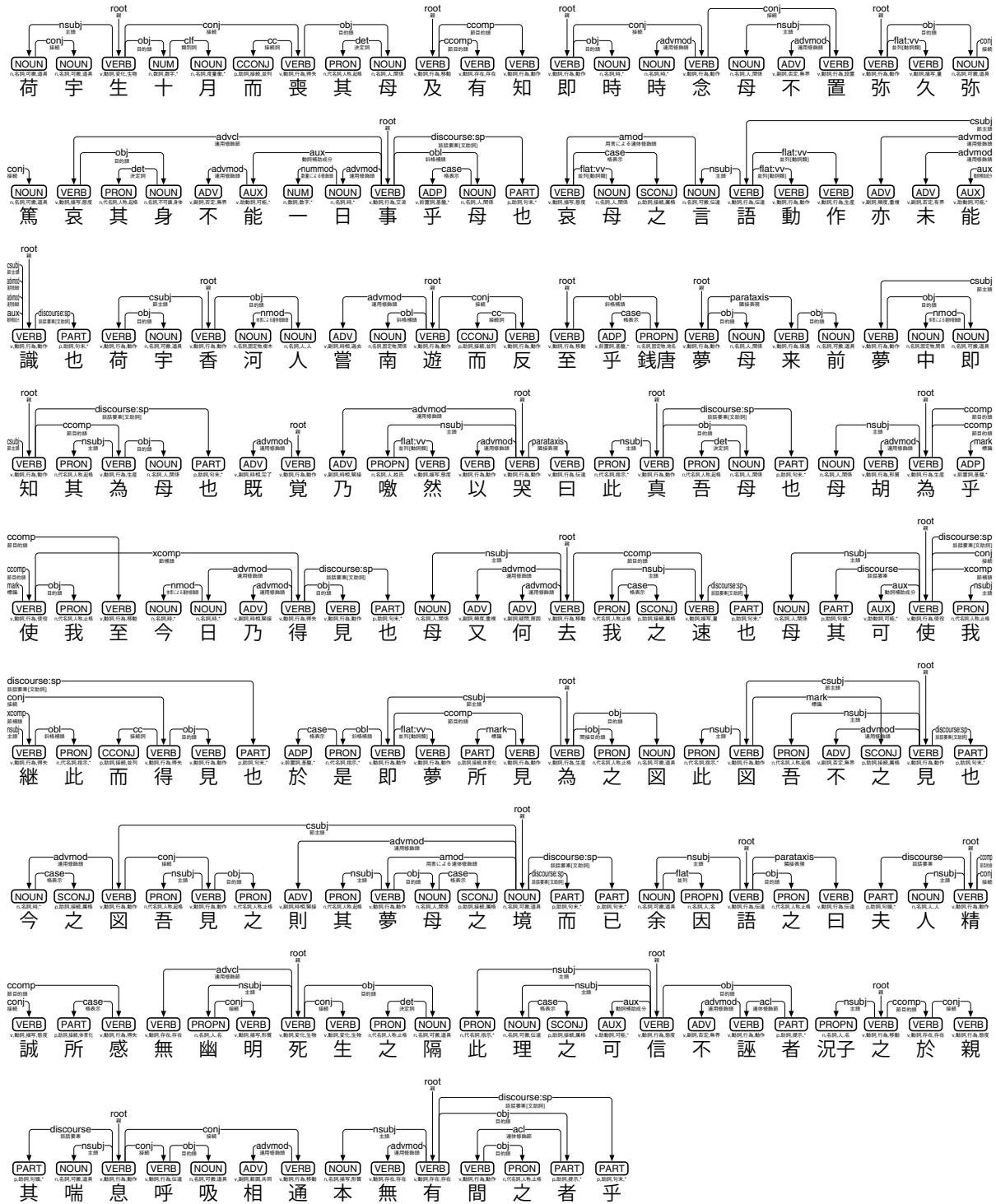


図 79: 手法Aの処理結果(2016年)

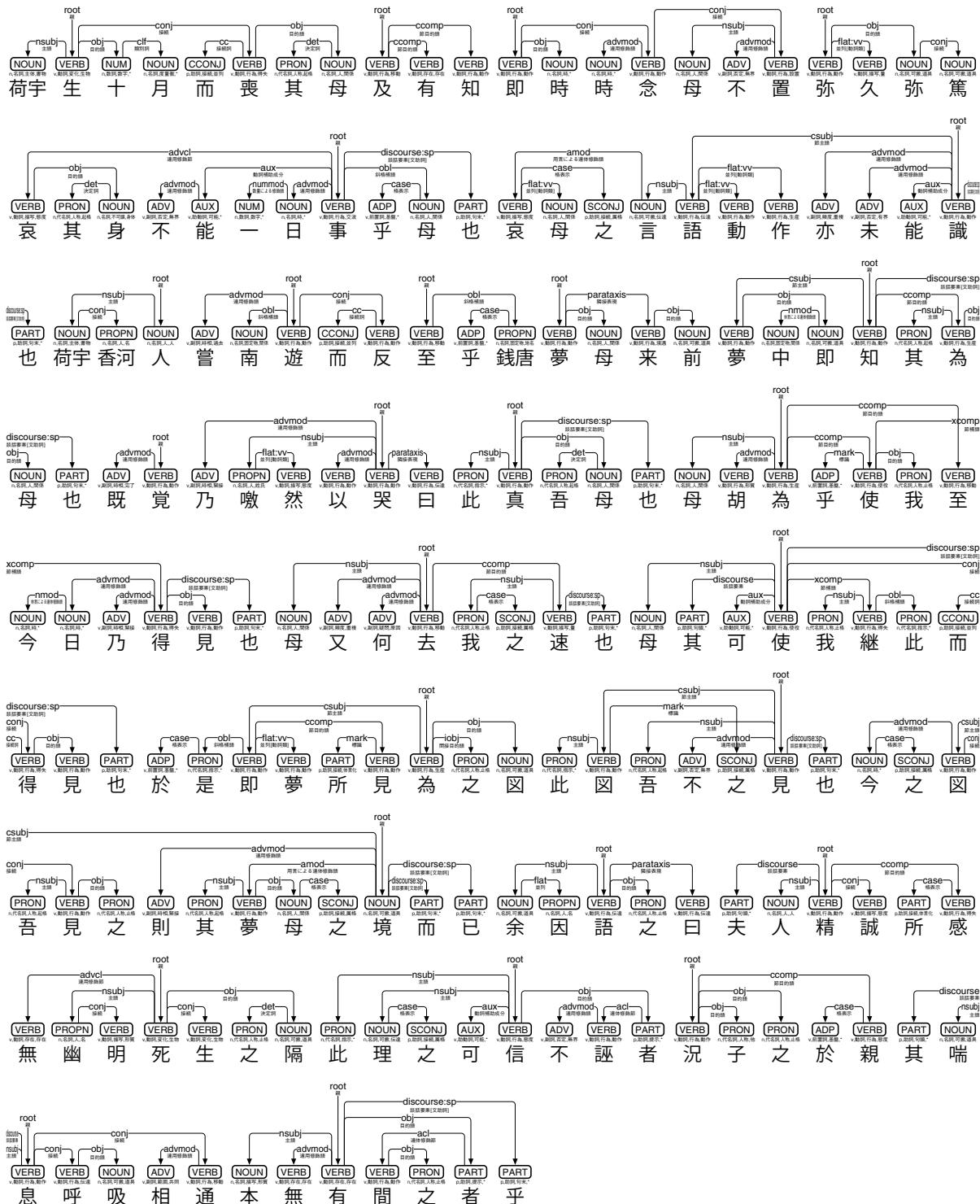


図 80: 手法Bの処理結果(2016年)

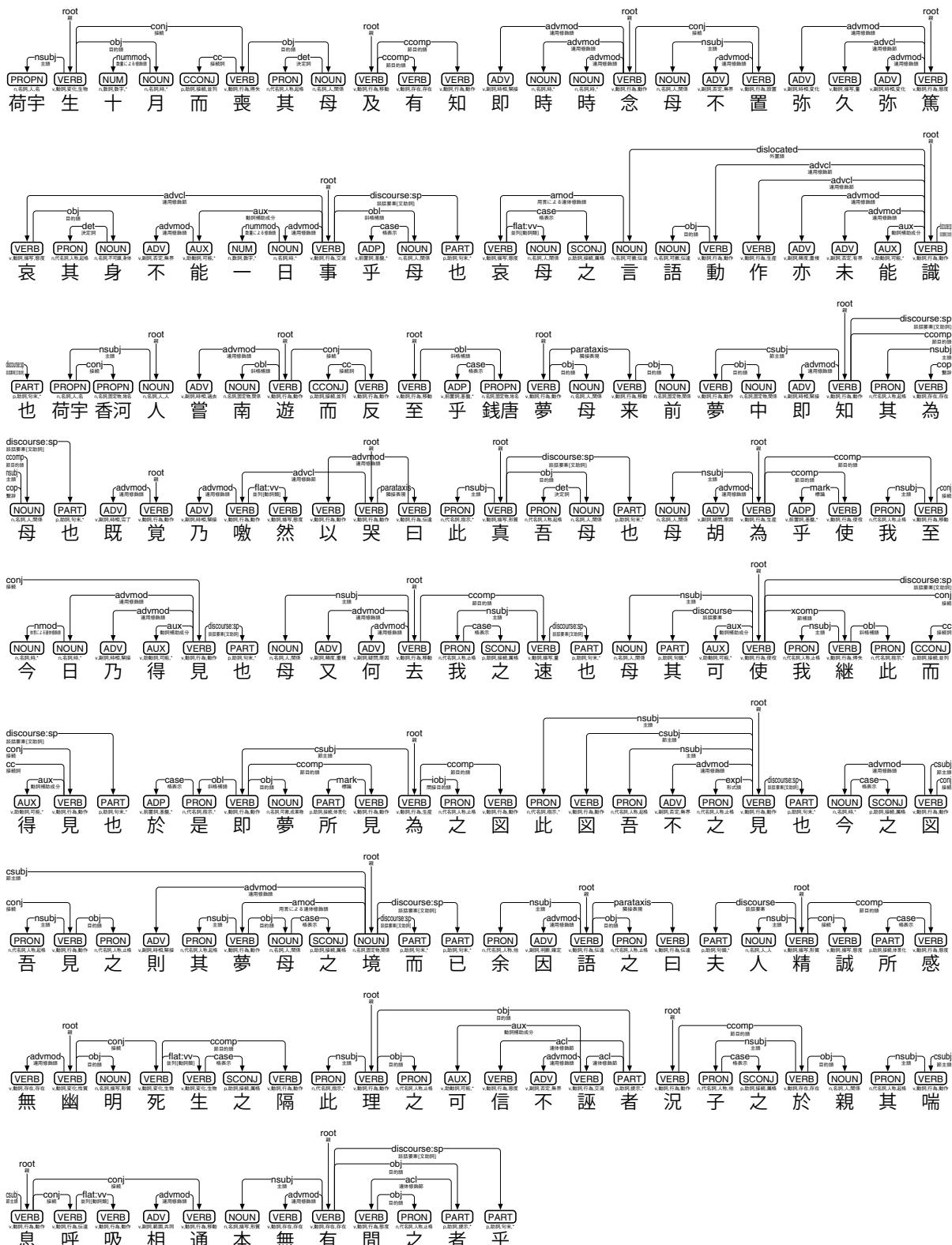


図 81: 手法Cの処理結果(2016年)

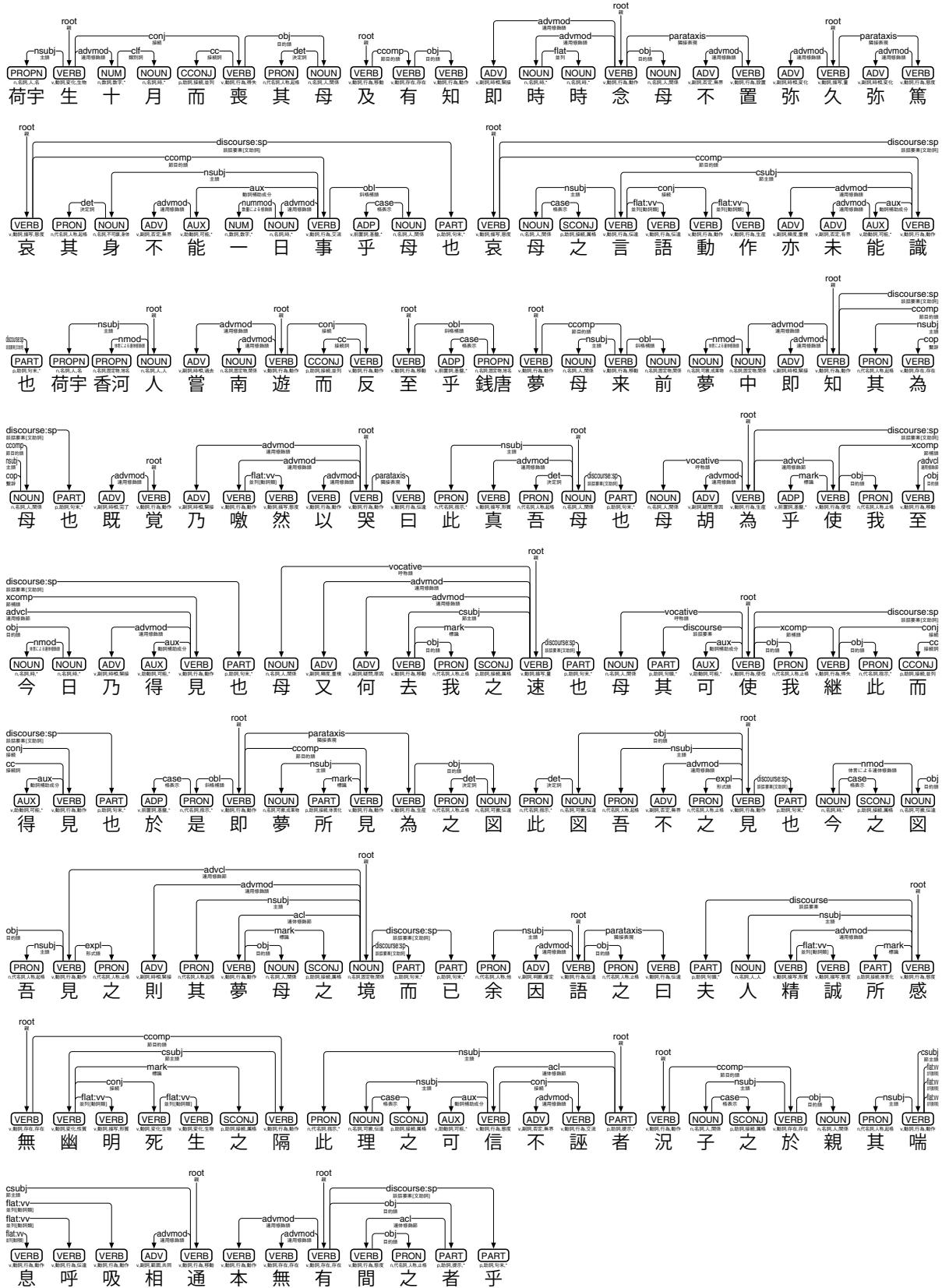


図 82: 手作業で作成した「正解」UD (2016 年)

假^は此^ニ異類^(e)已[。]
 則^チ世^ニ之^ニ為^ル人^ニ親^与子[、]而^有不^ニ慈^不孝^者、豈^獨愧^ニ于^古人[。]亦^タ

家^ニ蓄^ニ一^老狸^リ奴^レ將^レ誕^レ子^ヲ矣^(a)。一女童誤^{リテ}触^レ之^ニ而^{シテ}墮^ス日夕鳴^を
 鳴然[。]然[。]會^{タマタマ}有下餽^リ餽^ル兩^小狸^奴者[。]其始^ノ蓋^シ漠然^(注3)不^ニ相能^一也[。]老狸奴^{ナル}
 者[、]從^{ヒテ}而撫^レ之^ヲ、傍^は徨^う焉^(注4)躊躇^{ちよく}焉[。]臥^{スレバ}則^シ擁^レ之^ヲ、行^{ケバ}則^シ翊^レ之^ヲ。舐^ニ其^乳
 而讓^ル之^ニ食^一兩^小狸^奴者[、]亦久^シ而相忘^ル也[。]稍^シ即^レ之^ニ、遂^ニ承^レ其^乳
 焉[。]自^レ是[。]欣然^(注6)以^テ為^ス良^己之^母老狸奴^者、亦^タ居^(注7)然^{トシテ}以^テ為^ス良^己
 出^{ダスト}^(b)也[。]吁^{ああ}、亦^タ異^{ナル}哉^{かな}。

昔^ハ漢^明^(注8)德^馬后^ニ無^レ子[。]顯^宗取^リ他^人子[。]命^{ジテ}養^{ハシテ}之^ヲ曰^{ハク}人^子何^B
 必親生[。]但恨^ム愛^之不^レ至^耳^(C)。后^遂盡^{クシテ}心^撫育^シ而^{シテ}章^帝亦^タ恩^(注10)性[。]

天^至母^子慈^孝、始終^無纖^{せん}芥^{かい}之間[。]狸^奴之^事、適^有契^{かなフ}焉^(d)。然^{しかラバ}

図 83: 大学入試センター試験『国語』(2015年1月17日)第4問本文

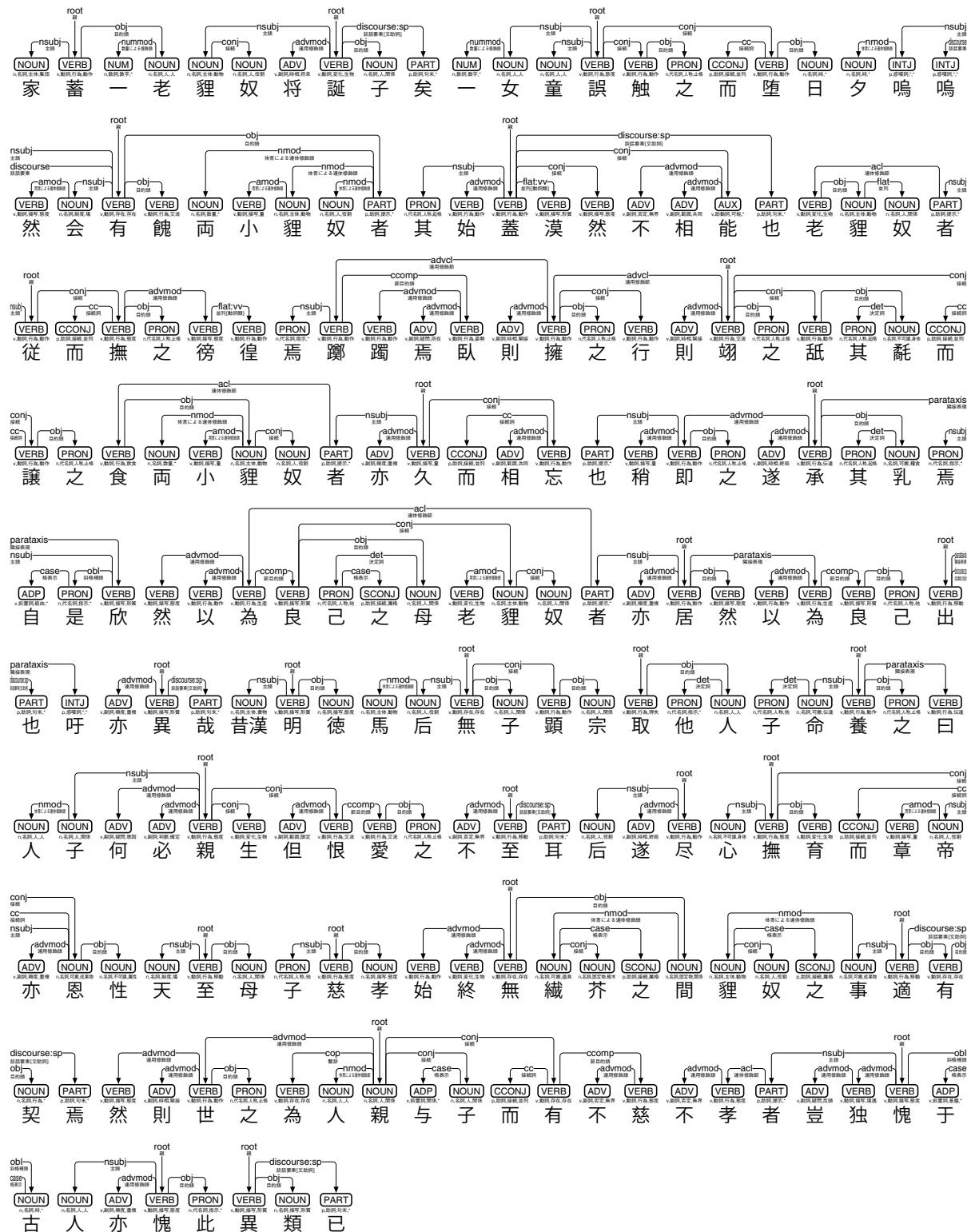


図 84: 手法①の処理結果(2015年)

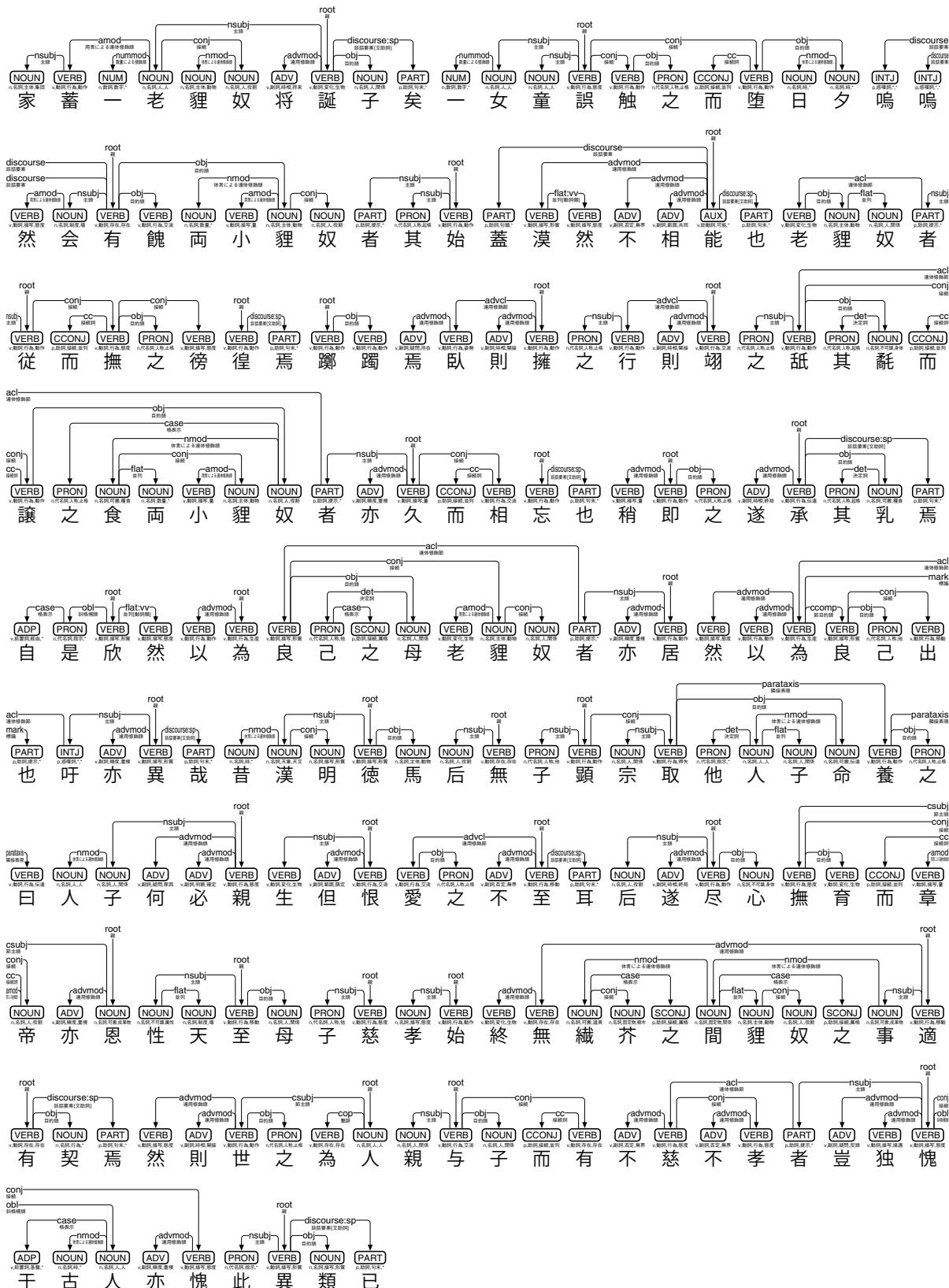


図 85: 手法②の処理結果(2015年)

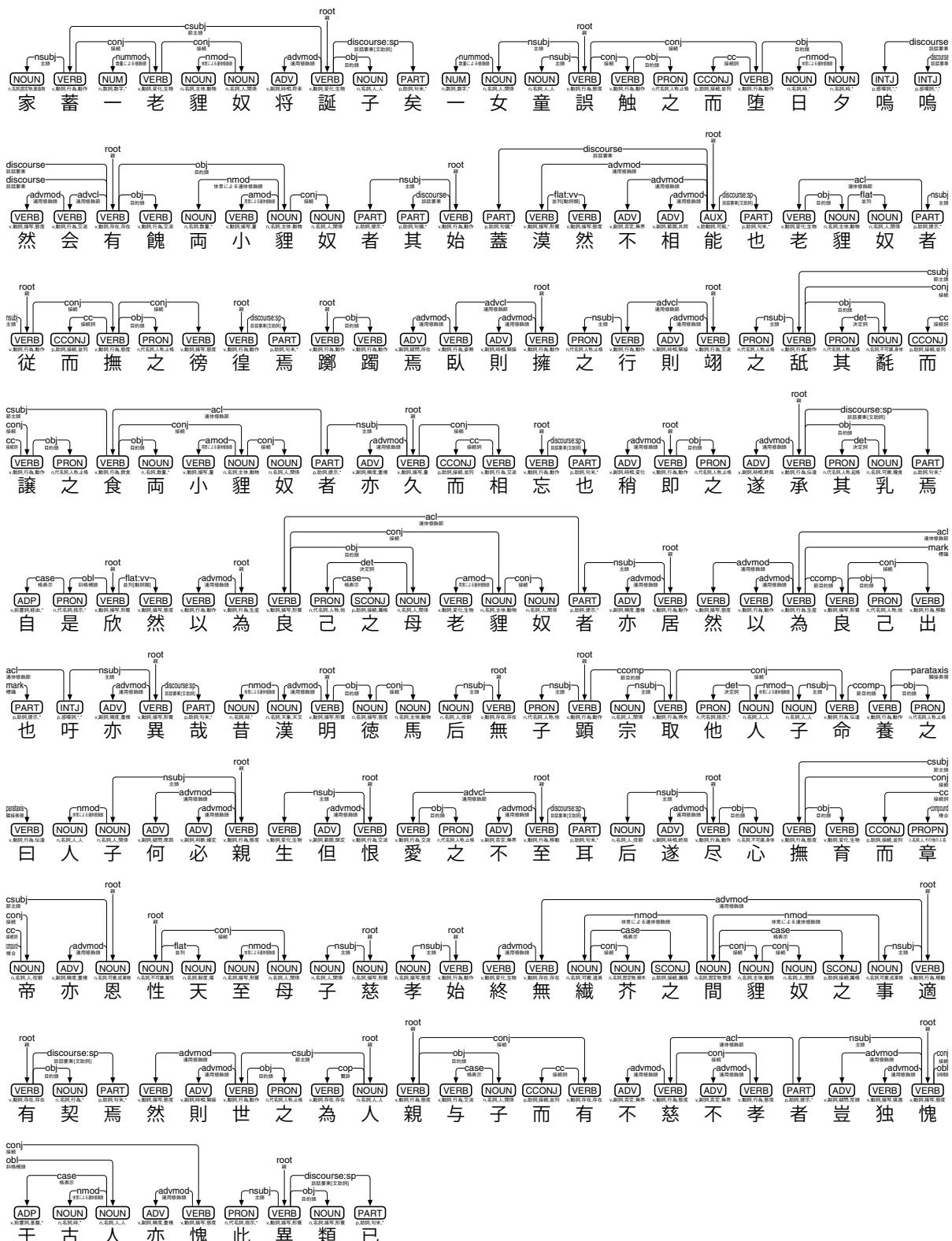


図 86: 手法③の処理結果 (2015 年)

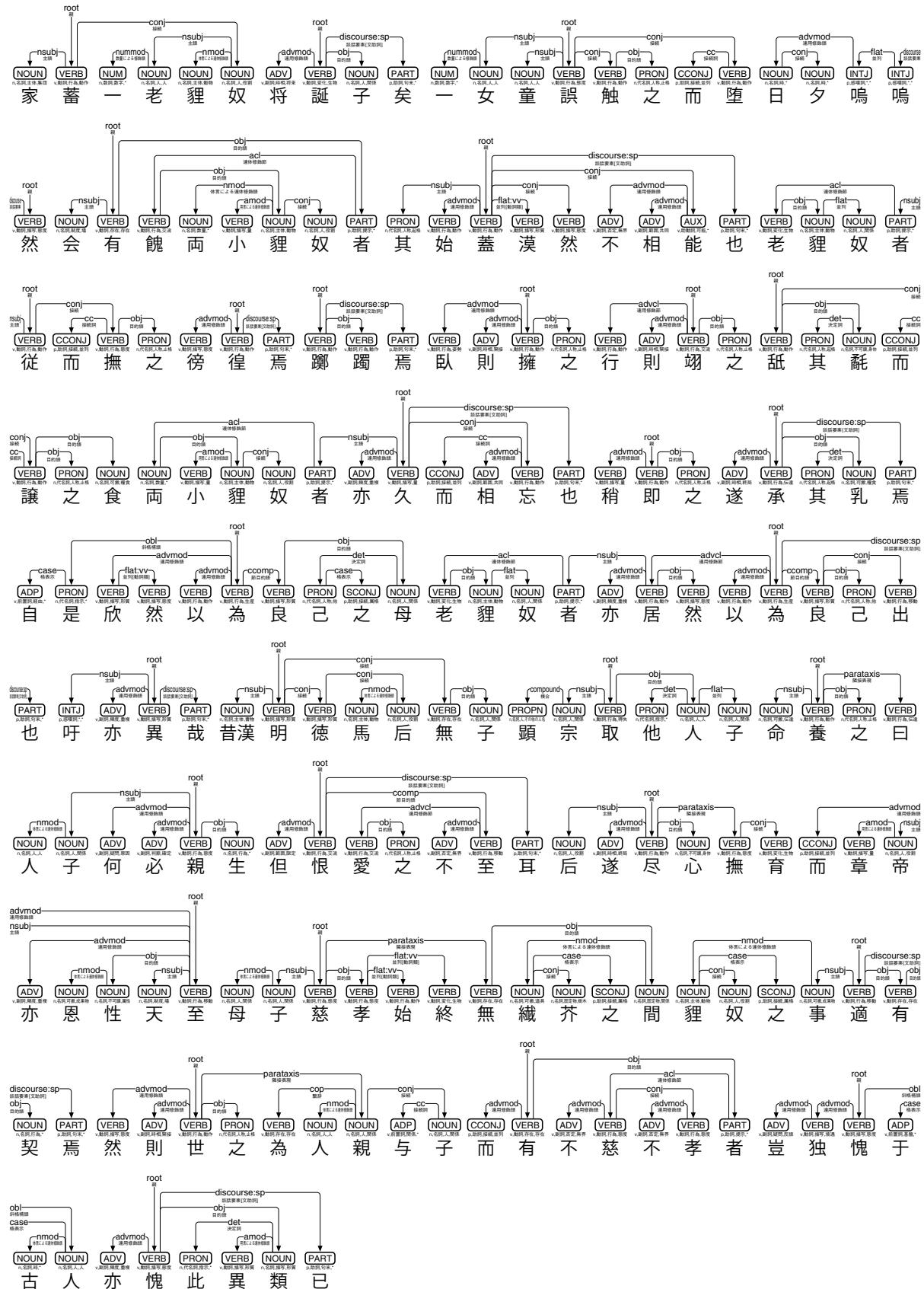


図 87: 手法①の処理結果(2015年)

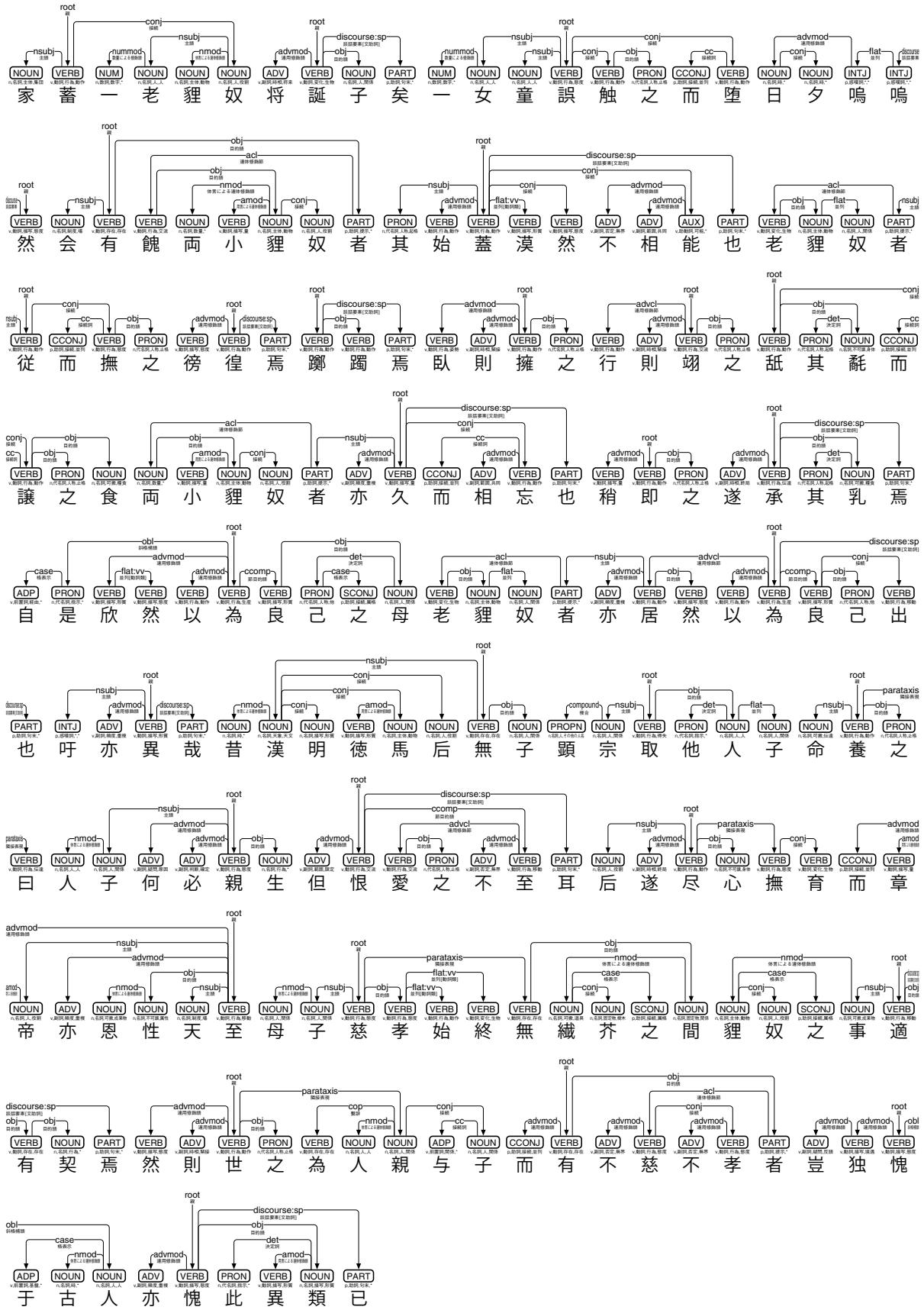


図 88: 手法②の処理結果(2015年)

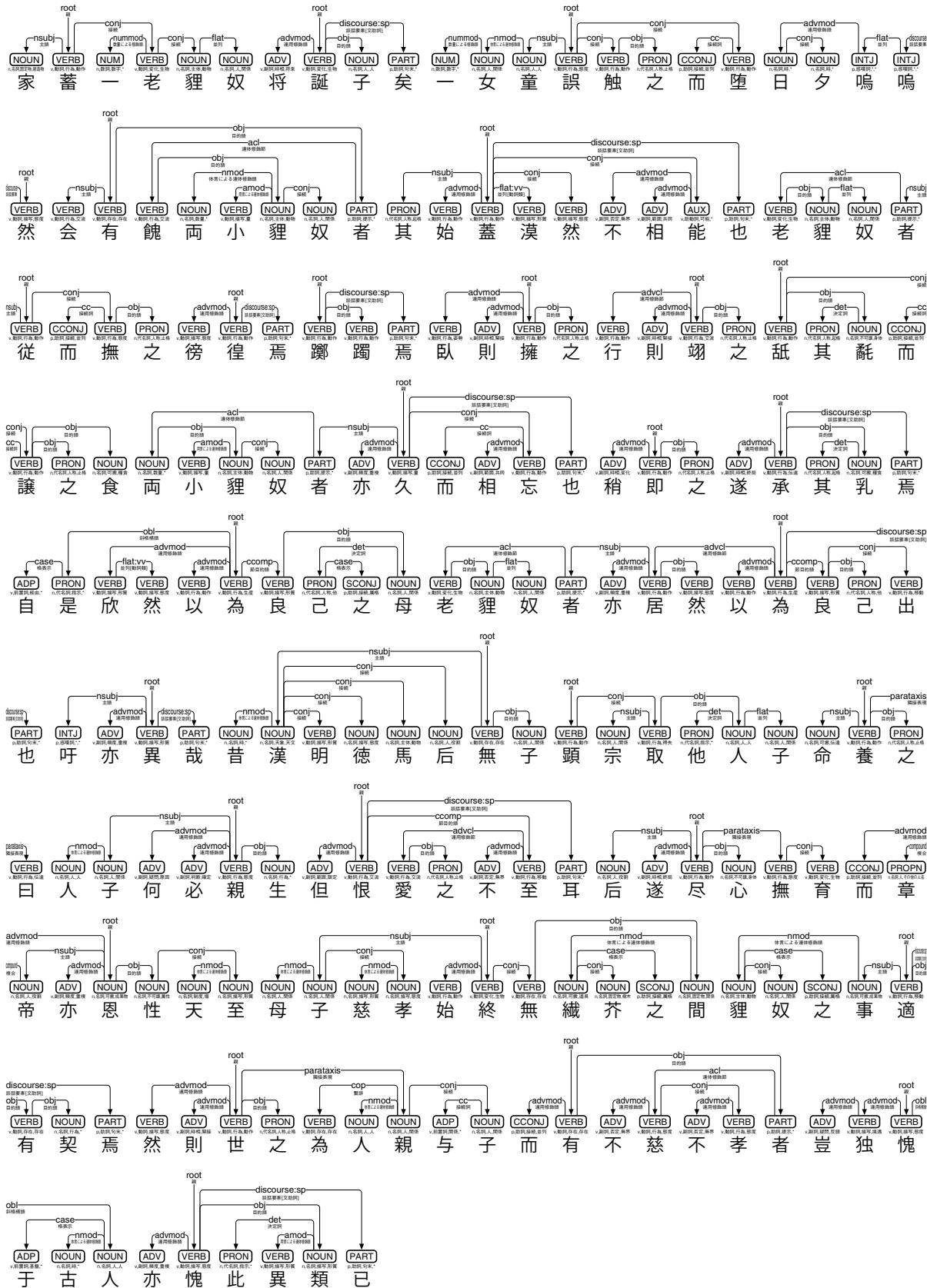


図 89: 手法③の処理結果(2015 年)

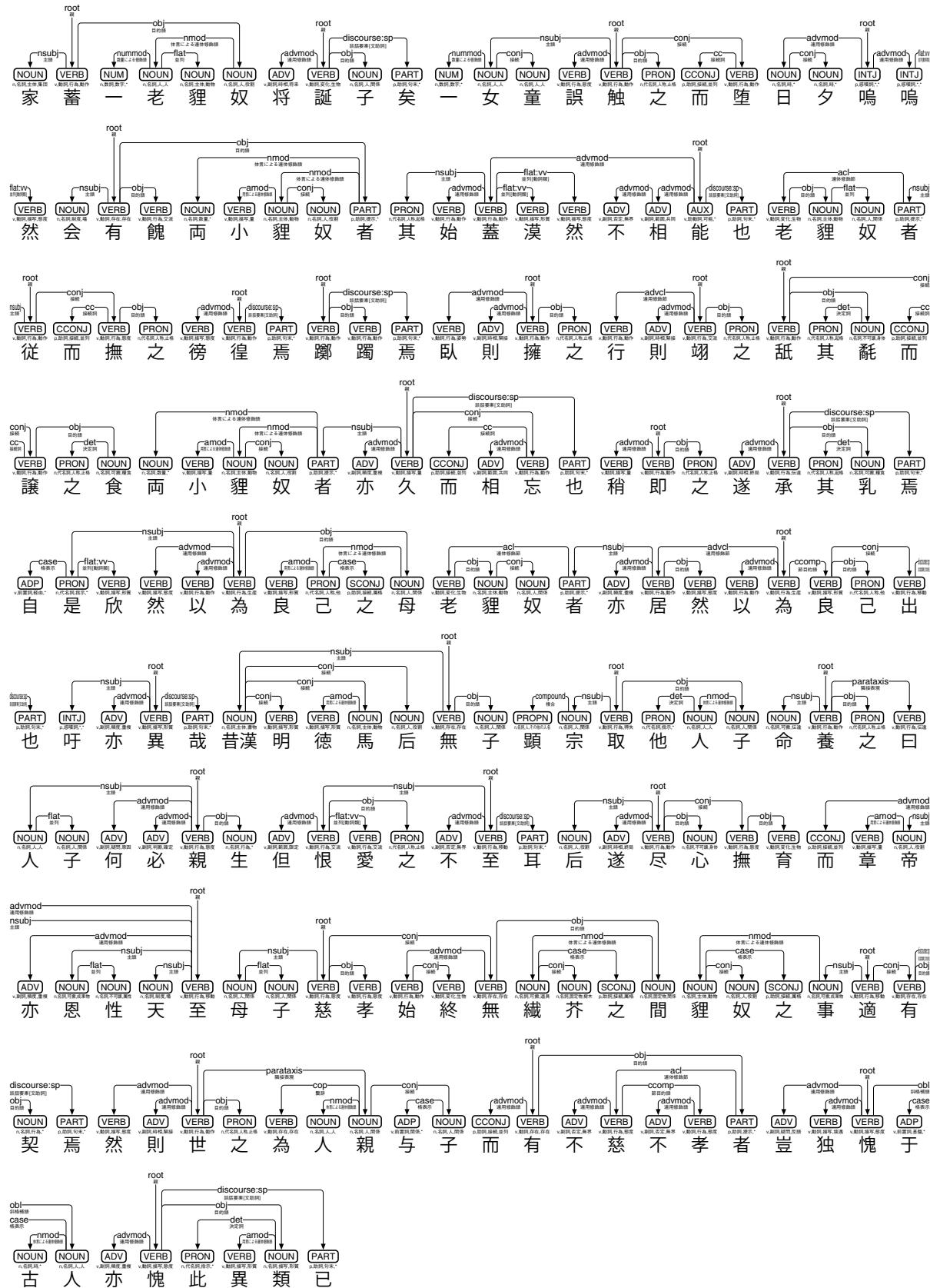


図 90: 手法①の処理結果(2015年)

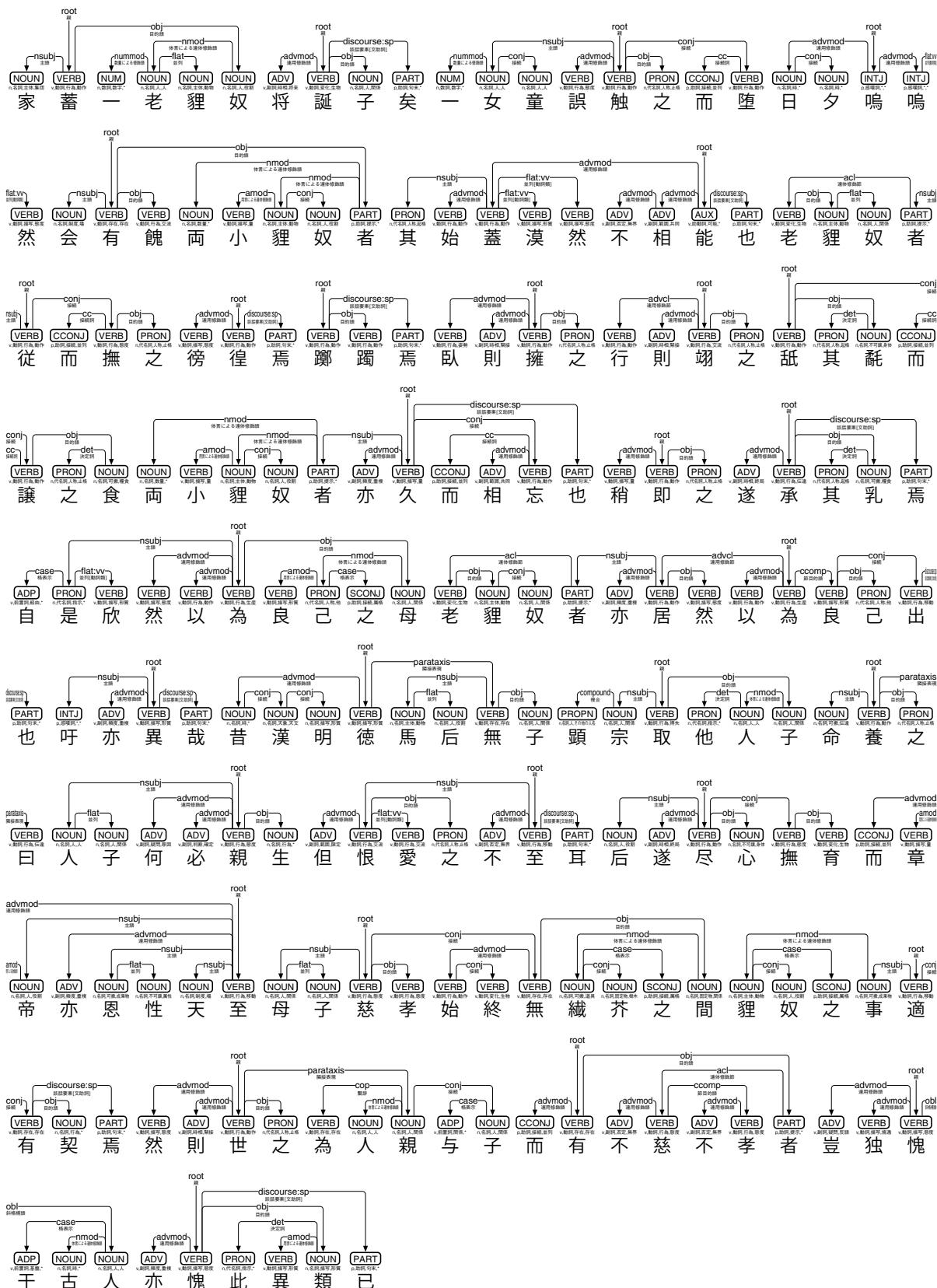


図 91: 手法②の処理結果(2015年)

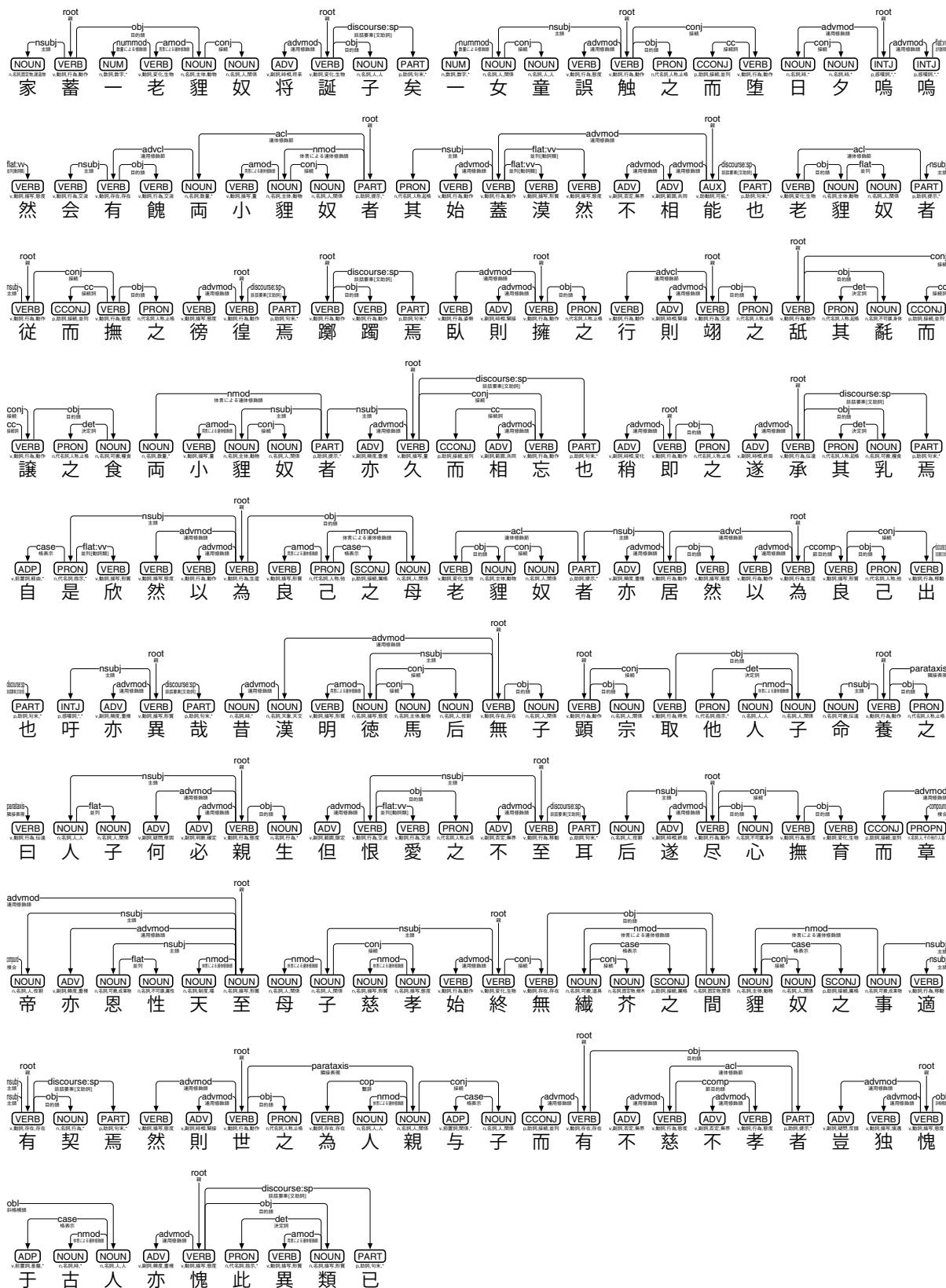


図 92: 手法③の処理結果(2015年)

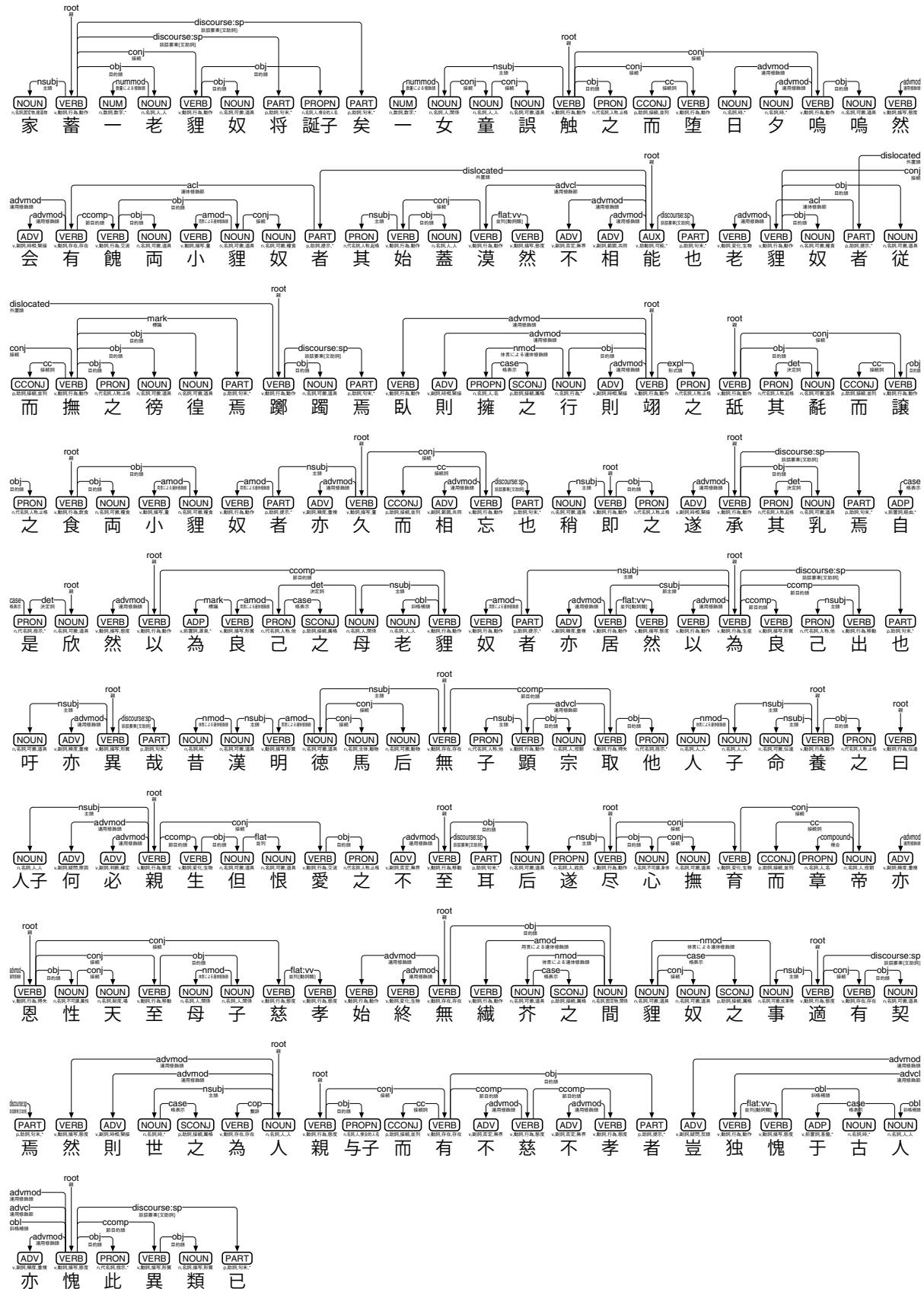


図 93: 手法Ⓐの処理結果(2015年)

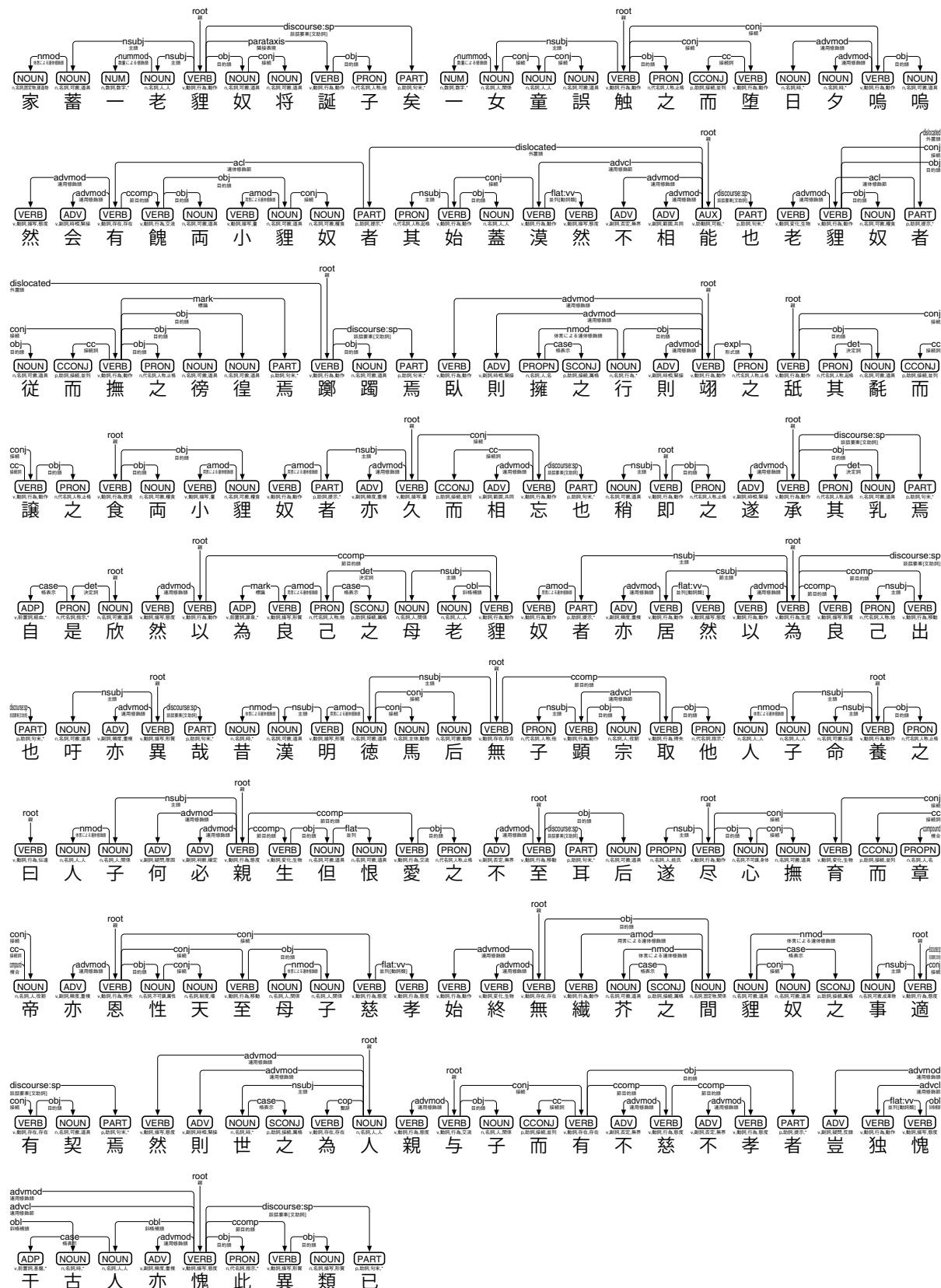


図 94: 手法②の処理結果(2015年)

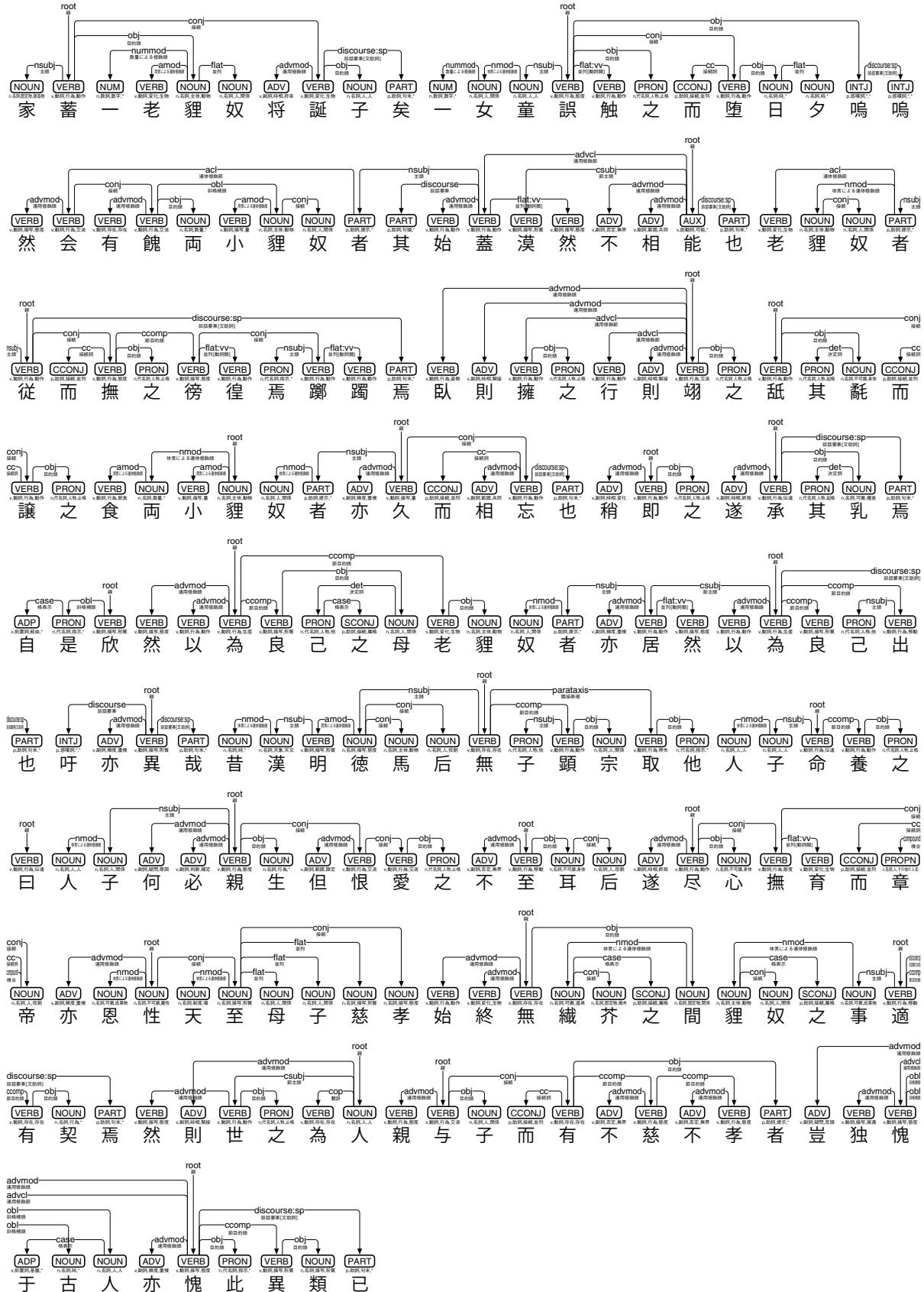


図 95: 手法④の処理結果(2015年)

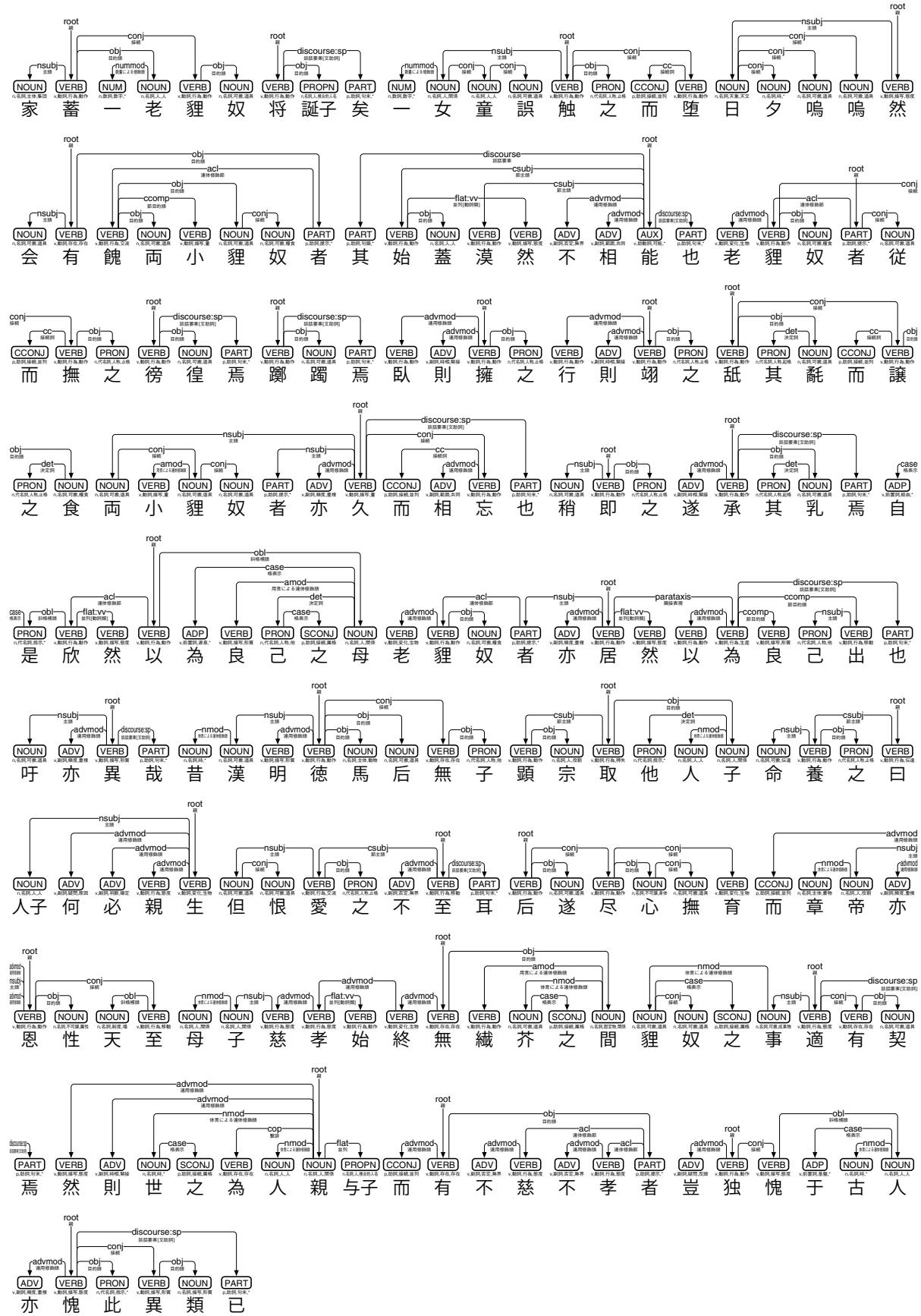


図 96: 手法Aの処理結果(2015年)

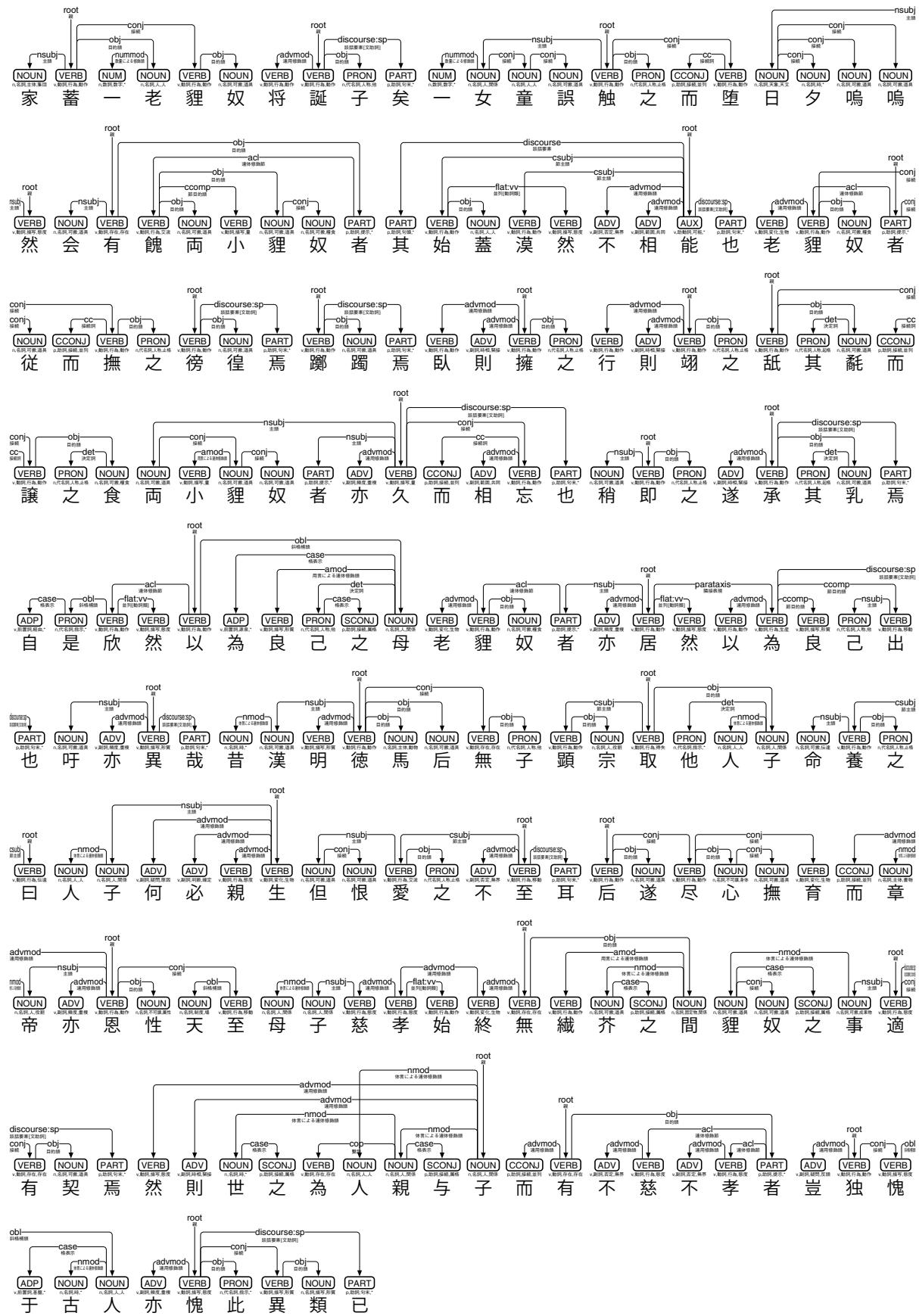


図 97: 手法Bの処理結果(2015年)

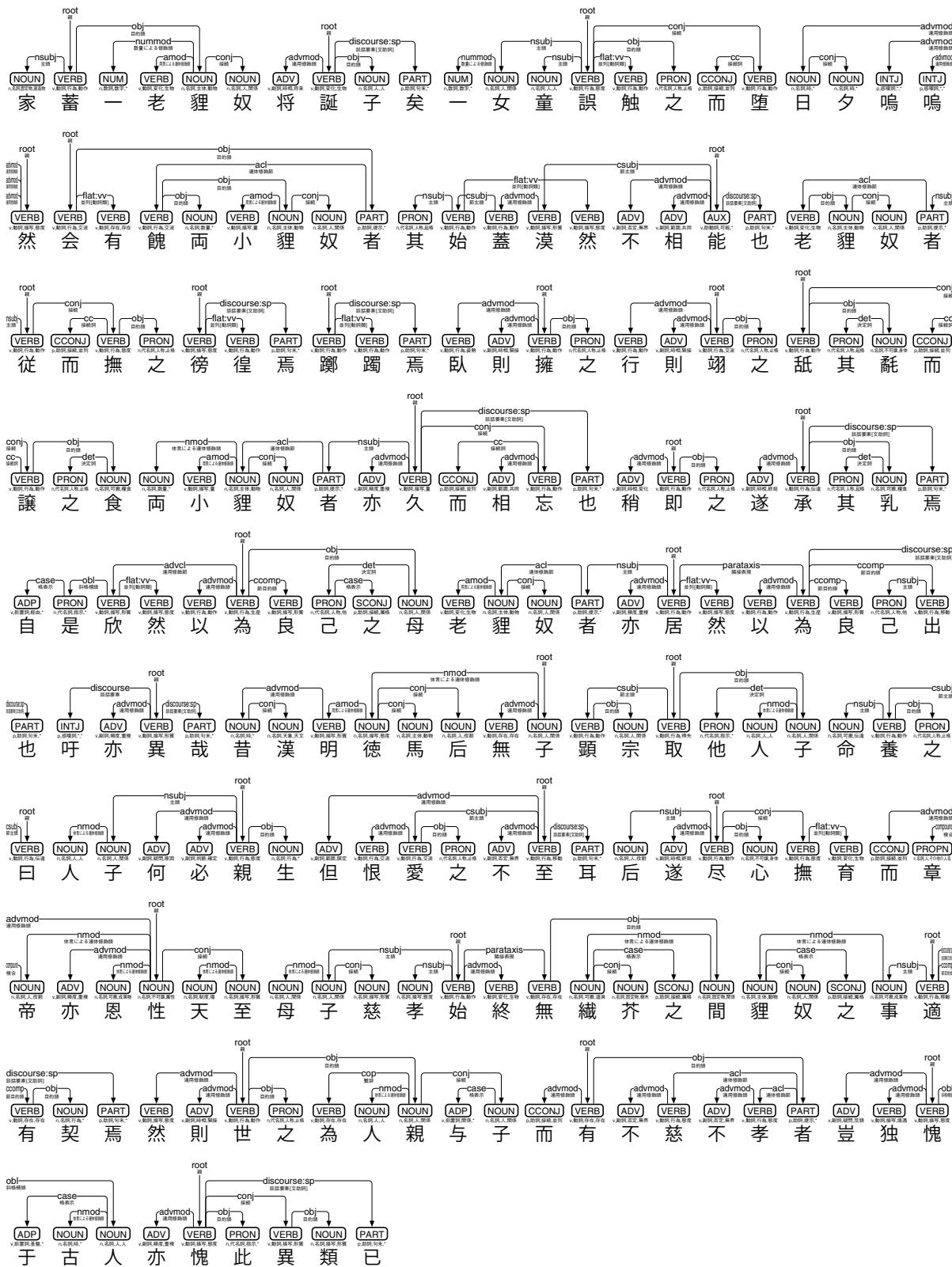


図 98: 手法Cの処理結果(2015年)

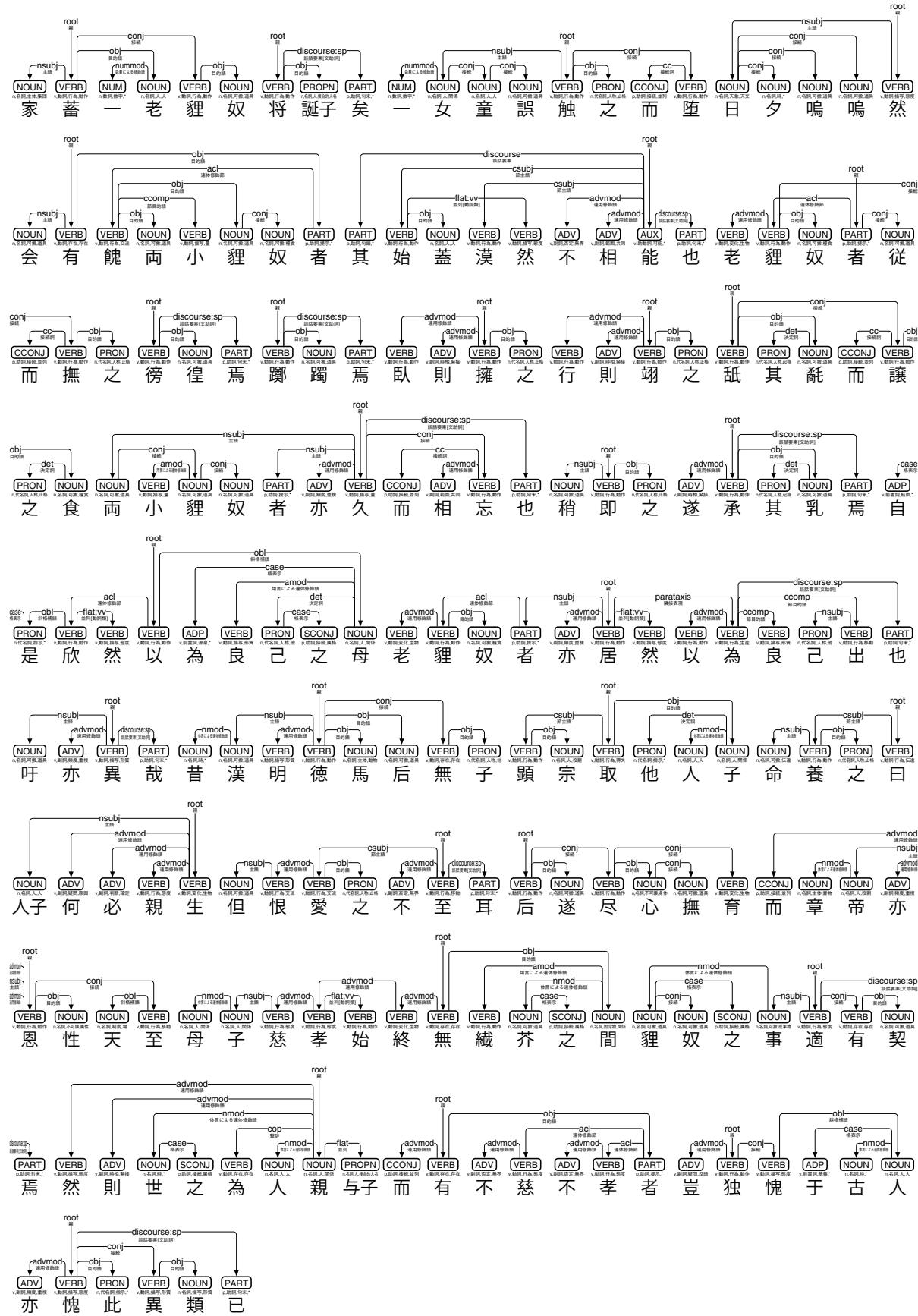


図 99: 手法Aの処理結果(2015年)

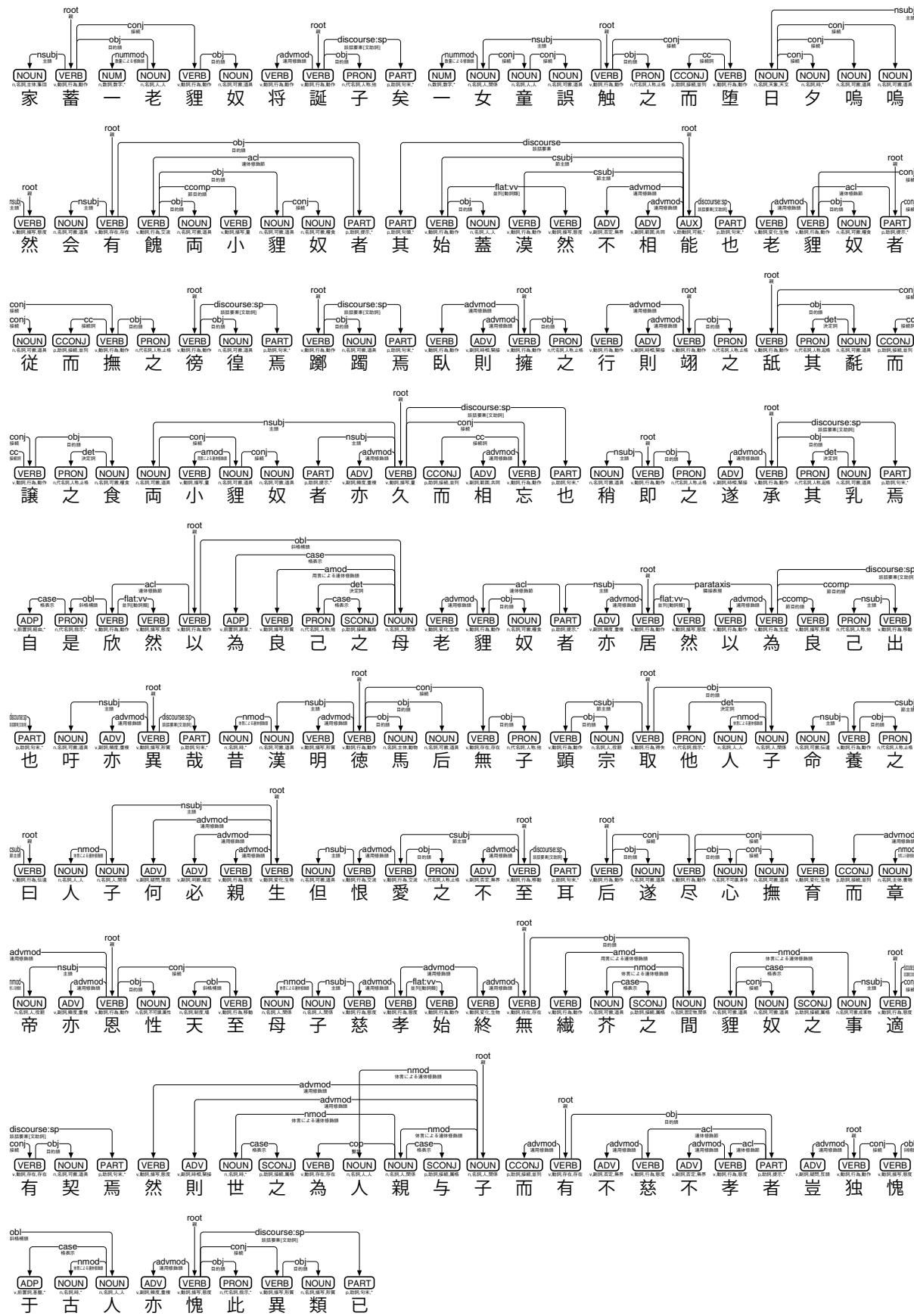


図 100: 手法Bの処理結果(2015年)

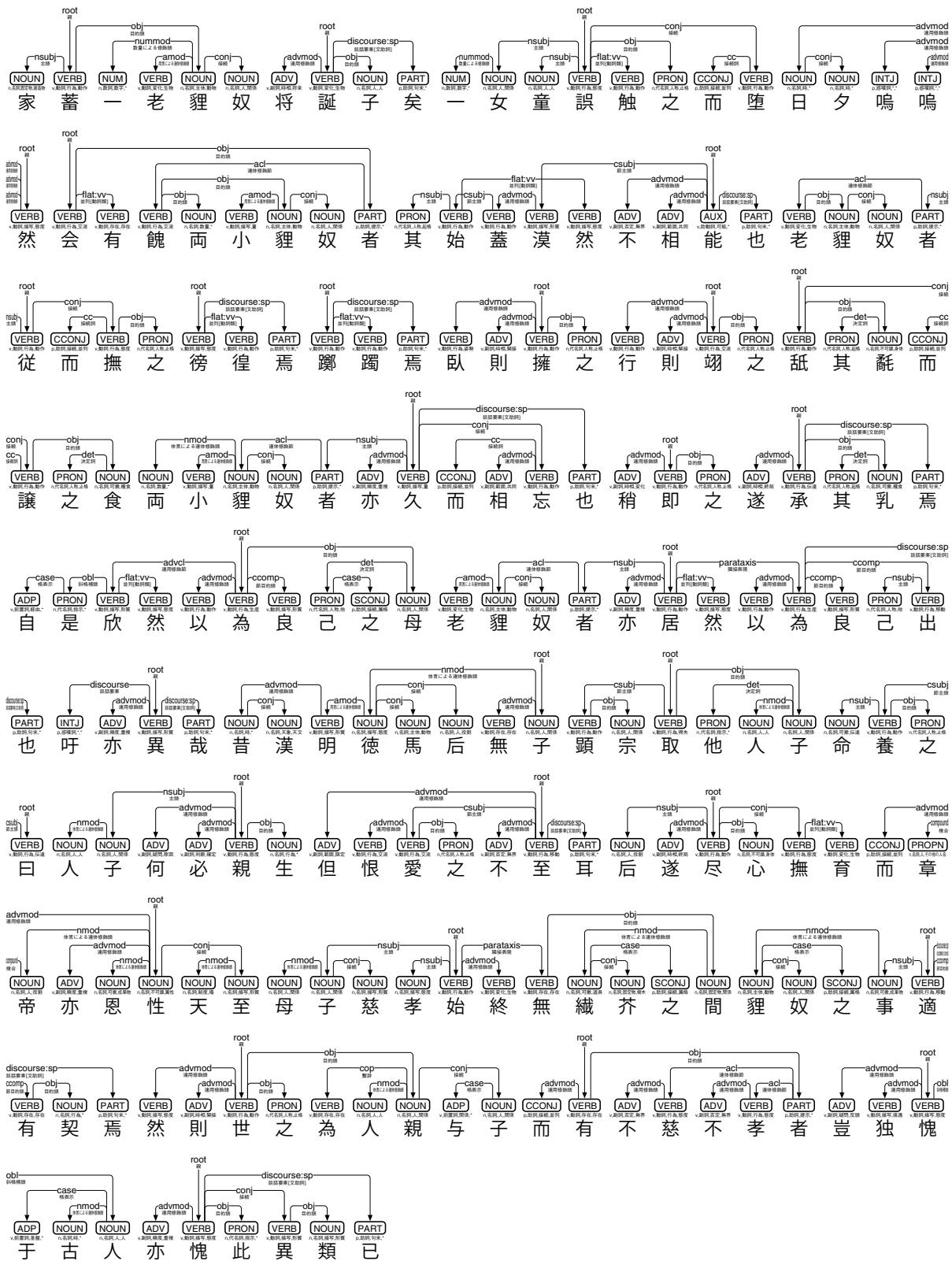


図 101: 手法Cの処理結果(2015年)

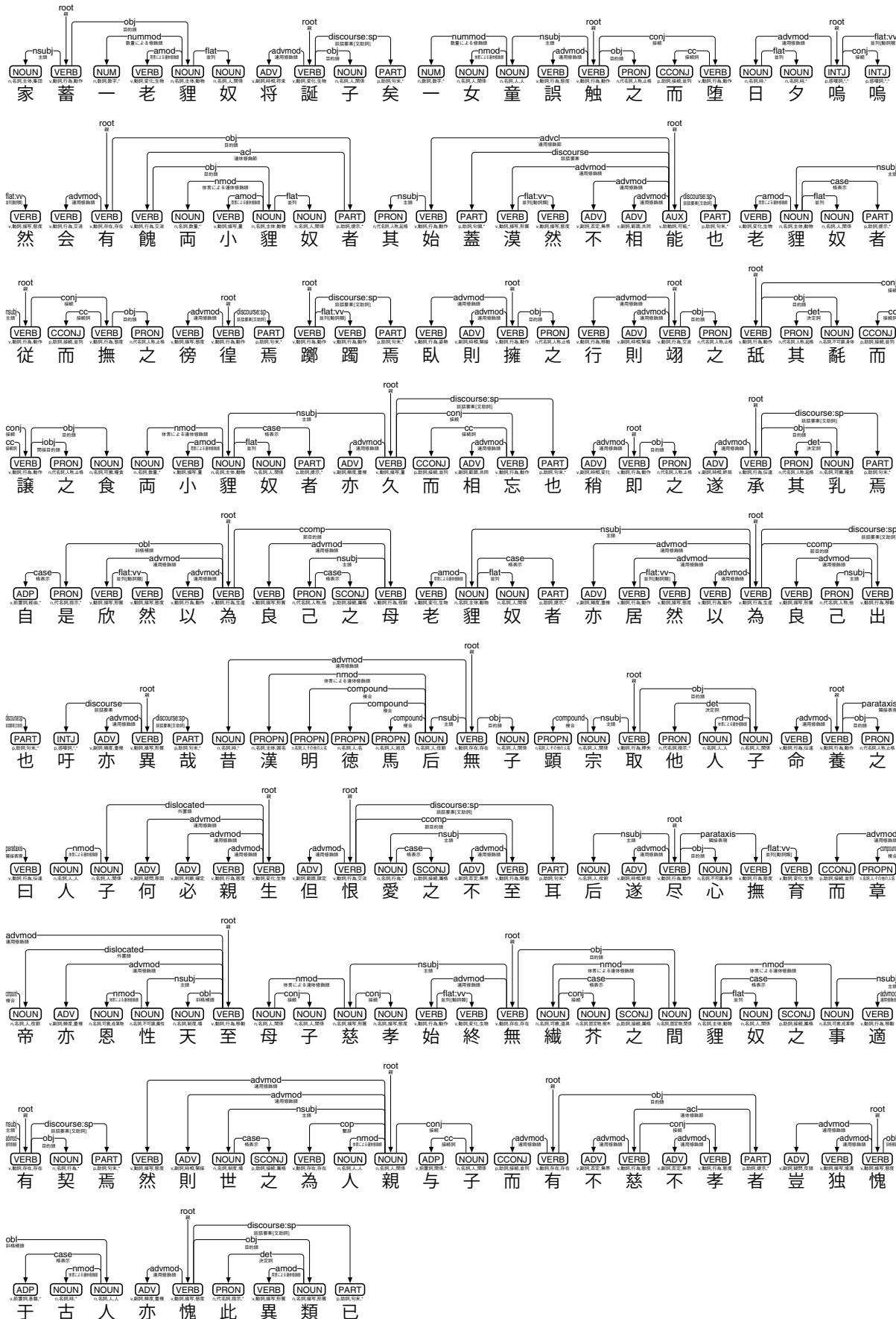


図 102: 手作業で作成した「正解」UD (2015 年)