



Deposited via The University of Sheffield.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/227925/>

Version: Published Version

---

**Article:**

Jiang, X. and Liu, J. (2023) Extracting the evolutionary backbone of scientific domains: The semantic main path network analysis approach based on citation context analysis. *Journal of the Association for Information Science and Technology*, 74 (5). pp. 546-569. ISSN: 2330-1635

<https://doi.org/10.1002/asi.24748>

---

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

# Extracting the evolutionary backbone of scientific domains: The semantic main path network analysis approach based on citation context analysis

Xiaorui Jiang<sup>1</sup>  | Junjun Liu<sup>2</sup>

<sup>1</sup>Research Centre for Computational Sciences and Mathematical Modelling, Coventry University, Coventry, UK

<sup>2</sup>Independent Researcher, Jiaxing, China

## Correspondence

Xiaorui Jiang, Coventry University, Coventry, UK  
Email: [xiaorui.jiang@coventry.ac.uk](mailto:xiaorui.jiang@coventry.ac.uk)

## Funding information

National Office for Philosophy and Social Sciences, Grant/Award Number: 18ZDA238

## Abstract

Main path analysis is a popular method for extracting the scientific backbone from the citation network of a research domain. Existing approaches ignored the semantic relationships between the citing and cited publications, resulting in several adverse issues, in terms of coherence of main paths and coverage of significant studies. This paper advocated the semantic main path network analysis approach to alleviate these issues based on citation function analysis. A wide variety of SciBERT-based deep learning models were designed for identifying citation functions. Semantic citation networks were built by either including important citations, for example, extension, motivation, usage and similarity, or excluding incidental citations like background and future work. Semantic main path network was built by merging the top- $K$  main paths extracted from various time slices of semantic citation network. In addition, a three-way framework was proposed for the quantitative evaluation of main path analysis results. Both qualitative and quantitative analysis on three research areas of computational linguistics demonstrated that, compared to semantics-agnostic counterparts, different types of semantic main path networks provide complementary views of scientific knowledge flows. Combining them together, we obtained a more precise and comprehensive picture of domain evolution and uncover more coherent development pathways between scientific ideas.

## 1 | INTRODUCTION

There were many methods to extract the evolutionary pathways between scientific ideas based on citation network analysis, such as algorithmic historiography (Garfield et al., 2003) and scientific historiograms (Lucio-Arias & Leydesdorff, 2008). Recently, main path analysis (MPA), originally proposed in Hummon and Doreian (1989), has become popular for extracting the major

knowledge diffusion paths among the main ideas advancing an analyzed scientific domain, since Batagelj (2003) proposed the efficient *search path counting* algorithms to weight citation edges and Verspagen (2007) laid out the algorithmic foundations for main path extraction.

Most MPA methods were citation semantics-agnostic, that is, ignoring the semantic relationships between publications. A direct consequence is semantically *incoherent* main path. Figure 1 illustrates a potential cause of this

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2023 The Authors. *Journal of the Association for Information Science and Technology* published by Wiley Periodicals LLC on behalf of Association for Information Science and Technology.

problem—inappropriate search path counts (SPC). In the top-right schematic image, the citation edges (A, B) and (B, C) are both background citations (“Neutral”) while the citation edge (A, C) is an extension citation (“Extends”). Ignoring citation function, we have  $SPC(A, B) \geq SPC(A, C)$  because the former is the sum of the number of paths through  $A \rightarrow B \rightarrow C$ , which is equal to  $SPC(A, C)$ , and the number of paths through  $A \rightarrow B \rightarrow X (X \neq C)$ . So traditional MPA approaches will select (A, B), but it is more reasonable to include the extension citation (A, C). Some studies adjusted citation weight by, for example, considering citation preferences according to discipline and publication time (Yu & Pan, 2021) or scaling search path count using citing publication’s prestige (Yu & Sheng, 2021). However, the problem was not solved. For example, if B is highly cited, then Yu and Pan’s approach will still choose (A, B) in main path exploration. Some weighing schemes used measures of similarity between the abstracts of citing and cited publications (Chen et al., 2022; Huang et al., 2022; Liu et al., 2014). However, such (indirectly inferred) similarity measures shall be less precise than authors’ own (directly stated) rationales to cite, aka *citation function* (Iqbal et al., 2021; Kunnath et al., 2022; Lyu et al., 2021).

Theoretically, traditional MPA approaches also tend to prefer long local paths.<sup>1</sup> Figure 1 illustrates this case. The left-most image shows a vanilla (semantics-agnostic) main path network (MPN). The longest local path from A00-2018 to D07-1096 is very stretched: distance (A00-2018, D07-1096) = 16. It is questionable whether knowledge indeed flows along such long paths with many unimportant citations such as “Neutral.” The middle image shows a snapshot of the semantic main path network (semantic MPN) extracted by considering extension (“Ext”) and motivation (“Mot”) citations. The path becomes more compact: distance (A00-2018, D07-1096) is decreased to 5. For another example, by further considering usage (“Use”) and similarity (“Sim”) citations, the longest distance from W96-0213 to W05-0516 is reduced from 17 to 5.

To the best of our knowledge, this is the first paper which marries citation function classification to MPA. We proposed a systematic approach to semantic main path network analysis (Section 4) based on citation function classification (Section 3), which solves both issues raised above. Multiple semantic citation networks were built using different citation functions, for which multiple semantic main path networks were extracted, assuming that different semantic networks capture different types of knowledge flows between different knowledge entities, such as ideational basis, methodological extension, tool usage, and similarity in problem or methodology, and so on. We conjecture that different semantic main path networks will collectively provide a more comprehensive representation of an analyzed domain. Note that, there were also some recent studies relying on citation importance classification (Ghosal et al., 2022; Hassan et al., 2018). Essentially, these approaches weighted citation edges by 1 (important) or 0 (incidental), screened out unimportant citations, did not further processing for knowledge flow analysis. The current paper is methodologically different. Citation function classification provides us with more flexible ways to perform MPA. The superiority of the proposed approach was qualitatively justified using two case studies (Section 5). In Section 6, this paper proposed a three-way quantitative evaluation framework. To the best of our knowledge, this is the first study about quantitative evaluation of MPA results. Experiments proved that extracting and merging multiple semantic main path networks achieved better (topical) coverage, (topical) coherence and (ranking) pertinence (Section 6).

## 2 | RELATED WORK

### 2.1 | Topological approaches of main path analysis

According to Verspagen (2007), MPA has two steps: citation weighting and main path extraction. Refer to Liu

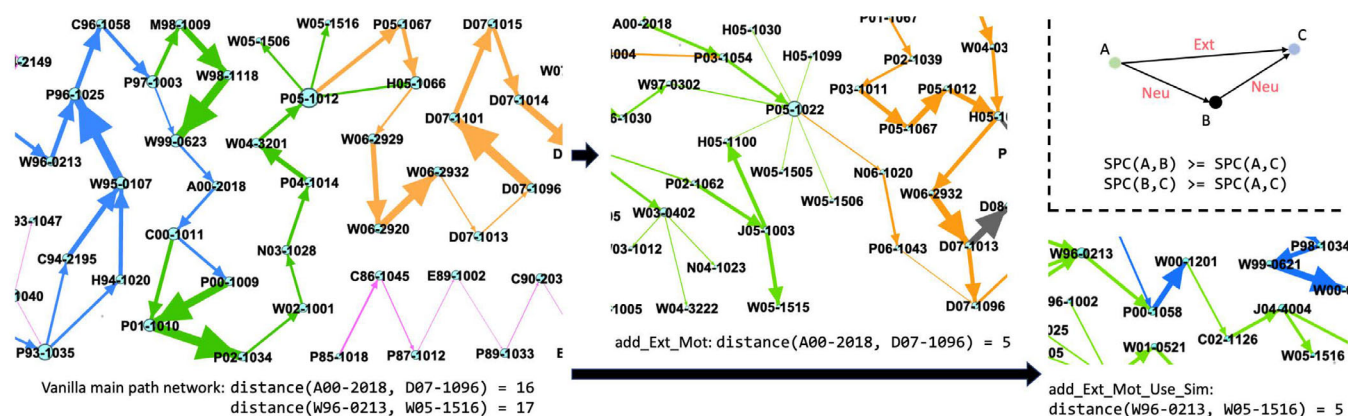


FIGURE 1 Motivations for semantic main path network analysis.

TABLE 1 Search path counting methods for main path analysis.

Method	Origins	Targets	Citation network standardization	
NPPC	All nodes	All	N/A	N/A
SPLC	All	Sinks (zero-outdegree)		Connect $s^*$ (resp. $t^*$ ) to all nodes (resp. sinks)
SPNP	All	All	Add a pseudo-source $s^*$ and a pseudo-sink $t^*$	Connect $s^*$ and $t^*$ to all nodes
SPC	Sources (zero-indegree)	Sinks		Connect $s^*$ (resp. $t^*$ ) to all origins (resp. sinks)

et al. (2019, 2020) for the discussions of best practices of each step. Citation weighting is traditionally based on each edge's traversal count in the search paths between a set of origin nodes and target nodes in a (usually reversed) citation network. We call them *topological* approaches. The ground-breaking work of Hummon and Doreian (1989) defined three measures: Node Pair Project Count (NPPC), Search Path Link Count (SPLC), and Search Path Node Pair (SPNP). SPLC is predominantly used today. Batagelj (2003) proposed an efficient unified algorithm based on "standardizing" citation networks (summarized in Table 1), and proposed the fourth measure Search Path Count (SPC). For each citation edge  $(u, v)$  in a standardized citation network, the citation weight is equal to the number of paths from pseudo-source to  $u$  multiplied by the number of paths from  $v$  to pseudo-sink. As citation networks are mostly acyclic, the calculation is done iteratively based on topological sort. Kuan (2020) empirically discussed the choices of these weighting variants. Several adjustments exist. Liu and Kuan (2016) proposed to decay search path by length with the belief that knowledge diffusion has higher information loss along long paths, while Yu and Sheng (2021) used citing papers' citation influence for adjustment.

Typically, main path extraction starts from certain chosen startpoints and greedily searches the highest weighted citation edges to follow. Verspagen (2007) enumerated paths from the source(s) with the maximal out-going edge weight as startpoint(s) so the main paths were called *forward local main paths* (Liu & Lu, 2012). Batagelj (2003) also tried the longest path as the *global main path* (Batagelj, 2003). Liu and Lu (2012) defined two new types of local main paths. *Backward local main path* starts from sinks and represents the significant knowledge flow from past to the most recent studies. They also found that these methods often miss the most significant citation edges, called *key-routes*, they proposed the fourth alternative called *key-route main path* which searches forward and backward simultaneously from key-routes. To increase the comprehensiveness of the extracted main paths, Liu and Lu (2012) heuristically selected the top- $K$  startpoints or key-routes and merged the main paths extracted from them. Recently, Chen et al. (2022) proposed a more efficient dynamic programming algorithm for exhaustive main path extraction.

## 2.2 | Semantic approaches in main path analysis

Liu et al. (2014) pioneered to use (expert-assigned) *citation relevancy* to adjust traversal count-based citation weighting. Of course, it be replaced by any semantic relatedness measure. For instance, Huang et al. (2022) claimed that using the weighted sum of the textual and structural similarities between cited and citing publications lead to better convergence, that is, different slices of main path correspond well to different phases of domain development. Topic modeling is another popular semantic approach. Kim et al. (2022) used Latent Dirichlet Allocation (LDA) to analyze topic diffusion along main paths. Kim et al. (2018) used the Citation Influence Model, an extended LDA model which also models the generation process of each citing publication's citation mixture (Dietz et al., 2007), to measure citation weights by topic similarity. Chen et al. (2022) calculated the Cosine similarity between citing and cited articles' topic distribution obtained by Latent Semantic Indexing (Deerwester et al., 1990). Notably, the citation relevancy of  $(u, v)$  is the sum of the pair-wise similarities between  $v$  and all other nodes  $u'$  on the current path toward  $v$ . While this treatment theoretically ensured the topical coherence of main path, it looks more straightforward to extract main paths from topic subnetworks and merge them. Community detection could be seen as an alternative way of finding topic subnetworks (Kim & Shin, 2018; Yu & Pan, 2021). To the best of our knowledge, citation function classification (Kunnath et al., 2022; Lyu et al., 2021) has never been applied to main path analysis before.

## 3 | CITATION FUNCTION CLASSIFICATION

### 3.1 | Dataset and annotation schemes

We created a large citation function dataset by merging and reannotating six existing datasets in the computational linguistics domain: Teufel2010 (Teufel, 2010; Teufel et al., 2006a), Dong2011 (Dong & Schäfer, 2011), Jha2016

(Abu-Jbara et al., 2013; Jha et al., 2017), Alvarez2017 (Hernández-Alvarez et al., 2017), Jurgens2018 (Jurgens et al., 2018), and Su2019 (Su et al., 2019). The source papers were crawled from ACL anthology.<sup>2</sup> Different annotation guidelines were adopted so all citation contexts were reannotated according to Teufel et al.'s 12-class annotation scheme (Teufel et al., 2006b) plus a “Future” class about future work. Reannotation is detailed in Supplementary Section B.1.<sup>3</sup> Some minority classes were still small, so we merged “PModi” with “PBas” into “Basis,” and reannotated “CoCo.” into “CoCoGM” or “CoCoRes.” This resulted in our own 11-class annotation scheme, which was also mapped to 7-class and 6-class schemes by category merging. Table 2 shows the statistics of our dataset Jiang2022.

### 3.2 | Citation function classification models

For the purpose of recognizing citation functions more correctly, a series of deep learning models were developed. SciBERT (Beltagy et al., 2019) was used to encode citation context, currently fixed to 2 and 3 sentences to each side of the citation sentence (*citance*). Three types of features were generated from the SciBERT-encoded context: (a) the *citation* representation  $\mathbf{h}$ , from the citation segment (represented by a pseudo-word “CITSEG”), (b) the *citance* representation<sup>4</sup>  $\mathbf{s}$ , pooled by citance encoder from the citation sentence, and (c) the *context* representation  $\mathbf{c}$ , pooled by context encoder from the whole context. The final feature vector  $\mathbf{f}$  was the concatenation of the three:  $\mathbf{f} = [\mathbf{h}; \mathbf{s}; \mathbf{c}]$ . Citation representation is mandatory because different citations in the same citance should have different feature representations, but citance and context representations were optional.

We tested two types of citation contexts. In a *sequential context*, no “[SEP]” (sequence separator) was inserted to separate context sentences. In this case, citance and context representations were directly pooled from citance tokens and context tokens respectively. Two options of citance encoder were tested: max-pooling and self-attention (Munkhdalai et al., 2016). In a *hierarchical context*, “[SEP]” symbols were inserted after each context sentence. Sentence representations were pooled using sentence pooler, for which “[SEP]” was used as the third option in addition to max-pooling and self-attention, and context representation was pooled indirectly from the representations of all context sentences. There were in total 34 model variants.<sup>5</sup> Due to the large GPU time required for training, we cherry-picked a subset of 11 relatively promising variants, shown in Table 3, based on initial experiments of all model variants with the 11-class

scheme. Section 4.1 will discuss how to pick the appropriate models to perform semantic MPA based on per-class performance analysis of different models.

## 4 | SEMANTIC MAIN PATH NETWORK ANALYSIS

### 4.1 | Model selection: Precision or recall

Per-class performance analysis showed that no single best model could beat others on all citation functions or on all annotation schemes (Tables S1–S3). Therefore, we needed to choose the most appropriate model as a binary classifier for each specific citation function. The most pertinent citation function for MPA should be extension (“Basis”/“Extends”) of cited work, and motivation (“Motivation”) by previous studies. Figures 2 and 3 show the performances of these two classes’ top models. The darker the color, the higher the performance. Although the best *extension* model was model 4 (seed = 5,171, “seed =” omitted hereafter) with the 6-class scheme, its recall was less competitive. Considering the small size of the extension class, for example, only 4.33% in our dataset, we decided to slightly *weigh recall over precision* (*recall-oriented*) and F1. The final choice had a good F1 and the highest recall, that is, model 11 (47,353, in solid red rectangle) trained with the 6-class scheme. Taking a similar recall-oriented approach, we chose model 7 (32,491) trained with the 6-class scheme as the “best” *motivation* model.

We hoped that semantic citation networks could capture as many important citations as possible such as usage according to Valenzuela et al. (2015) and similarity according to Lu et al. (2014). For *usage* citations, we also took a recall-oriented approach. According to Figure 4, we opted for model 7 (13,249) trained with the 11-class scheme which achieved the highest F1, because the recall of the chosen model was already high enough and its precision was much higher than other candidates. To further enrich the semantic citation network, we decided to add *similarity* citations because Teufel’s annotation guidelines say similarity is between problems and solutions rather than results (Teufel, 2010). According to Figure 5, the selected model was model 11 (25,603) trained with the 11-class scheme.

The other way is to delete unimportant citations, for example, neutral citations (“Neutral”/“Background”) or future work citations (“Future”) in our case. Due to the dominant size of neutral citations and high performance on this class (Figure 6), we decided to *trade recall for precision* (*precision-oriented*) for neutral (“Neutral”/“Background”), so model 2 (5,171) with the 7-class scheme was

**TABLE 2** Statistics of the reannotated dataset Jiang2022 and citation function scheme mapping.

Original reannotations (12+1 class)				Our 11-class scheme <sup>a</sup>			Mapped to 7-class scheme <sup>b</sup>			Mapped to 6-class scheme		
Label	Size		Ratio citseg	Label	Size citseg	Ratio	Label	Size citseg	Ratio	Label	Size citseg	Ratio
	citstr	citseg <sup>c</sup>										
Future	97	85	2.21%	Future	85	2.21%	Future	85	2.21%	Future	85	2.21%
CoCoXY	200	152	3.94%	CoCoXY	152	3.94%	Background	1,773	46.00%	Background	1,615	41.90%
Neut	1,924	1,463	37.96%	Neutral	1,463	37.96%						
Weak	223	158	4.10%	Weakness	158	4.10%	ComOrCon	479	12.43%	ComOrCon	944	24.49%
CoCoGM	390	299	7.76%	CoCoGM	328	8.51%						
CoCo <sup>d</sup>	108	80	2.08%	CoCoRes	151	3.92%						
CoCoR0	107	100	2.59%				Similar	307	7.97%			
PSup	123	100	2.59%	Support	100	2.59%						
PSim	247	207	5.37%	Similar	207	5.37%	Motivation	288	7.47%	Motivation	288	7.47%
PMot	365	288	7.47%	Motivation	288	7.47%	Uses	755	19.59%	Uses	755	19.59%
PUse	794	755	19.59%	Usage	755	19.59%						
PModi	72	65	1.69%	Basis	167	4.33%	Extends	167	4.33%	Extends	167	4.33%
PBas	134	102	2.65%									
Total	4,784	3,854			3,854			3,854			3,854	

<sup>a</sup>CoCoXY is the Contrast/Comparison between two cited publications; CoCoGM/Res is the Comparison/Contrast between cited and citing publications Goals or Methods/Results; Basis is the Cited publication is ideationally based on; Support is the Cited and citing publications support each other's claims or can be computationally plugged into each other.

<sup>b</sup>ComOrCon is the Comparison/Contrast between citing and cited publications.

<sup>c</sup>A citseg (citation segment) is a number of consecutive citstrs (citation string) cited in the same place. Citation function classification is done for each citseg.

<sup>d</sup>“CoCo-” samples were re-annotated into either CoCoGM or CoCoRes based on what is compared.

TABLE 3 Selected citation function classification models.

ID	citation_encoder (h)	context_type	sentence_pooler	citance encoder (s)	context_encoder (c)
1	O (used)	Sequential	N/A	max_pool	max_pool
2-3	O	Sequential	N/A	X (not used)	max_pool (2); self_attend (3)
4-6	O	Sequential	N/A	max_pool (4); self_attend (5); X (6)	X
7-8	O	Hierarchical	max_pool	X	max_pool (7); self_attend (8)
9-11	O	Hierarchical	N/A	max_pool (9); self_attend (10); X (11)	X

Extends		Jiang2021 (11-class)					Jiang2021 (7-class)					Jurgens2018 (6-class)				
ID	metric	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353
4	precision	70.37	65.52	63.33	62.50	74.07	65.52	76.00	63.64	70.97	64.71	80.77	56.76	52.94	62.50	68.97
	recall	55.88	55.88	55.88	58.82	58.82	55.88	55.88	41.18	64.71	64.71	61.76	61.76	52.94	58.82	58.82
	f1 score	62.30	60.32	59.38	60.61	65.57	60.32	64.41	50.00	67.69	64.71	70.00	59.15	52.94	60.61	63.49
8	precision	70.82	43.75	61.29	80.95	71.43	57.89	75.00	75.00	75.50	63.33	50.00	61.29	62.96	75.00	60.00
	recall	50.00	41.18	55.88	50.00	44.12	64.71	61.76	52.94	61.76	55.88	58.82	55.88	50.00	44.12	52.94
	f1 score	58.62	42.42	58.46	61.82	54.55	61.11	67.74	62.07	67.74	59.38	54.05	58.46	55.74	55.56	56.25
11	precision	53.85	61.29	70.83	55.88	63.33	72.00	64.71	61.09	57.58	57.89	65.52	50.00	72.00	53.85	67.65
	recall	61.76	55.88	50.00	55.88	55.88	52.94	64.71	55.88	55.88	64.71	55.88	41.18	52.94	61.76	67.65
	f1 score	57.53	58.46	58.62	55.88	59.38	61.02	64.71	58.46	56.72	61.11	60.32	45.16	61.02	57.53	67.65

FIGURE 2 Performances of selected models for extension citations.

Motivation		Jiang2022 (11-class)					Jiang2022 (7-class)					Jurgens2018 (6-class)				
ID	metric	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353
6	precision	57.14	60.27	63.79	56.72	60.66	59.32	67.21	62.90	71.15	56.06	62.96	58.33	60.94	55.07	53.52
	recall	68.97	75.86	63.79	65.52	63.79	60.34	70.69	67.24	63.79	63.79	56.82	60.34	67.24	65.52	65.52
	f1 score	62.50	67.18	63.79	60.80	62.18	59.83	68.91	65.00	67.27	59.68	60.71	59.32	63.93	59.84	58.91
7	precision	60.27	60.00	55.56	59.68	58.73	64.71	66.67	62.96	62.50	55.74	56.42	54.22	57.35	68.25	60.78
	recall	75.86	62.07	60.34	63.79	63.79	56.90	65.52	58.62	68.97	59.62	67.24	77.59	67.24	74.14	53.45
	f1 score	67.18	61.02	57.85	61.67	61.16	60.55	66.09	60.71	65.57	57.14	61.42	63.83	61.90	71.07	56.88
11	precision	62.90	60.00	66.67	68.85	54.41	66.67	67.27	58.57	57.14	58.90	54.41	66.13	69.84	61.29	74.55
	recall	67.24	62.07	68.97	72.41	63.79	68.97	63.79	70.69	68.97	74.14	63.79	70.69	75.86	65.52	70.69
	f1 score	65.00	61.02	67.80	70.59	58.73	67.80	65.49	64.06	62.50	65.65	58.73	68.33	72.73	63.33	72.57

FIGURE 3 Performances of selected models for motivation citations.

selected. Because both precision and recall were high for future work citations (Figure 7), it was OK to adhere to the precision-oriented approach and select model 8 (32,941) with the 11-class scheme because it achieved high enough precision and the best F1.

## 4.2 | Semantic main path network extraction

### 4.2.1 | Citation network building

Starting from an empty citation network, a citation edge was added between a pair of publications if there existed

at least one in-text citation about extension **or** motivation (add\_Ext\_Mot) using the “best” *extension* or *motivation* models selected in the recall-oriented approach in Section 3.1. Taking the same recall-oriented approach, more citation edges were added if there existed at least one *usage* citation (plus\_add\_Use), and the semantic citation network was further expanded with *similarity* citations (plus\_add\_Sim). On the other hand, we also built the fourth semantic citation network by deleting unimportant in-text citations from the original citation network. For each pair of publications, if *all* in-text citations between them were *neutral* or *future work* citations, the citation edge was removed from the citation network (del\_Bkg\_Fut).

Usage ID	Jiang2022 (11-class)			Jiang2022 (7-class)			Jurgens2018 (6-class)									
	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353	
1	metric	73.65	73.33	74.48	75.84	71.52	74.50	75.71	78.29	76.80	73.20	74.65	79.07	76.06	80.95	78.69
	precision	72.19	80.13	71.52	74.83	74.83	73.51	70.20	66.89	63.58	74.17	70.20	67.55	71.52	56.29	63.58
	recall	72.91	76.58	72.97	75.33	73.14	74.00	72.85	72.14	69.57	73.68	72.35	72.86	73.72	66.41	70.33
2	precision	83.05	80.77	78.87	80.45	79.53	79.43	79.66	78.69	72.37	77.94	78.17	68.94	67.60	74.29	83.19
	recall	64.90	69.54	74.17	70.86	66.89	74.17	62.25	63.58	72.85	70.20	73.51	73.51	80.13	68.87	65.56
	f1 score	72.86	74.73	76.45	75.35	72.66	76.71	69.89	70.33	72.61	73.87	75.77	71.15	73.33	71.48	73.33
4	precision	79.39	79.67	77.08	74.64	78.17	78.99	75.54	70.55	73.10	79.85	82.03	81.54	72.79	83.74	76.67
	recall	68.87	64.90	73.51	68.21	73.51	72.19	69.54	76.16	70.20	70.86	69.54	70.20	70.86	68.21	76.16
	f1 score	73.76	71.53	75.25	71.28	75.77	75.43	72.41	73.25	71.62	75.09	75.27	75.44	71.81	75.18	76.41
7	precision	77.30	76.77	73.79	75.52	81.68	74.50	75.71	78.29	76.80	73.20	76.81	77.44	72.03	79.84	76.92
	recall	71.19	78.81	70.86	71.52	70.86	73.51	70.20	66.89	63.58	74.17	70.20	68.21	68.21	65.56	66.23
	f1 score	74.66	77.78	72.30	73.47	75.89	74.00	72.85	72.14	69.57	73.68	73.36	72.54	70.07	72.00	71.17

FIGURE 4 Performances of selected models for usage citations.

Similar ID	Jiang2022 (11-class)			Jiang2022 (7-class)			Jurgens2018 (7-class)				
	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353	
2	metric	57.45	46.94	65.79	60.00	67.50	62.07	53.12	54.29	59.65	63.46
	precision	64.29	54.76	59.52	71.43	64.29	58.06	54.84	61.29	54.84	53.23
	recall	60.67	50.55	62.50	65.22	65.85	60.00	53.97	57.58	57.14	57.89
5	precision	65.12	54.90	63.89	61.36	58.00	63.46	63.27	69.77	57.14	57.63
	recall	66.67	66.67	54.76	64.29	69.05	53.23	50.00	48.39	51.61	54.84
	f1 score	65.88	60.22	53.97	62.79	63.04	57.89	55.86	57.14	54.24	56.20
6	precision	58.14	56.82	58.54	62.79	60.53	59.32	70.00	60.71	66.67	81.08
	recall	59.52	59.52	57.14	64.29	54.76	56.45	45.16	54.84	58.06	48.39
	f1 score	58.82	58.14	57.83	63.53	57.50	57.85	54.90	57.63	62.07	60.61
11	precision	50.91	53.19	61.22	59.52	56.82	64.29	50.75	55.93	57.14	56.67
	recall	66.67	59.52	71.43	59.52	59.52	58.06	54.84	53.23	45.16	54.84
	f1 score	57.73	56.18	65.93	59.52	58.14	61.02	52.71	54.55	50.45	55.74

FIGURE 5 Performances of selected models for similarity citations.

Background ID	Jiang2022 (11-class)			Jiang2022 (7-class)			Jurgens2018 (6-class)									
	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353						
2	metric	75.50	76.14	76.04	78.85	75.43	82.17	78.02	80.72	78.59	78.06	76.74	75.08	74.34	74.34	
	precision	77.82	68.60	74.74	69.97	75.43	82.87	79.21	81.74	82.30	78.37	76.85	78.40	74.38	77.78	
	recall	76.64	72.17	75.39	74.14	75.43	82.52	80.23	79.84	81.50	78.48	77.45	77.56	74.73	76.02	
6	precision	76.36	71.78	72.78	78.75	75.77	82.07	82.91	80.11	79.94	76.07	78.22	76.38	76.88	73.08	76.56
	recall	71.67	79.86	78.50	77.13	75.77	75.84	81.74	83.71	80.62	84.83	78.70	76.85	75.93	82.10	75.62
	f1 score	73.94	75.61	75.53	77.93	75.77	78.83	82.32	81.87	80.28	80.21	78.46	76.62	76.40	77.33	76.09
8	precision	75.60	73.61	75.95	78.07	69.64	79.53	76.49	78.17	80.50	76.62	76.05	76.98	72.30	76.88	71.43
	recall	75.09	72.35	75.43	71.67	79.86	86.24	86.80	81.46	81.18	82.87	72.53	69.14	76.54	75.93	78.70
	f1 score	75.34	72.98	75.68	74.73	74.40	82.75	81.32	79.78	80.84	79.62	74.25	72.85	74.36	76.40	74.89
11	precision	74.38	73.11	76.22	73.75	75.68	77.15	78.98	82.20	77.40	81.65	67.94	76.20	75.00	76.19	75.37
	recall	82.35	76.11	79.86	75.77	76.45	88.20	78.09	77.81	83.71	75.00	82.41	78.09	80.56	74.07	79.32
	f1 score	78.12	74.58	78.00	74.75	76.06	82.31	78.53	79.94	80.43	78.18	74.48	77.13	77.68	75.12	77.29

FIGURE 6 Performances of selected models for neutral/background citations.

Future		Jiang2022 (11-class)					Jiang2022 (7-class)					Jurgens2018 (6-class)				
ID	metric	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353	5171	13249	25603	32491	47353
2	precision	87.50	84.62	66.67	76.47	72.22	75.00	70.59	92.31	72.22	76.47	75.00	100.00	87.50	81.25	76.47
	recall	82.35	64.71	94.12	76.47	76.47	70.59	70.59	70.59	76.47	76.47	70.59	64.71	82.35	76.47	76.47
	f1 score	84.85	73.33	78.05	76.47	74.29	72.73	70.59	80.00	74.29	76.47	72.73	78.57	84.85	78.79	76.47
5	precision	82.35	88.24	68.42	73.68	81.25	86.67	83.33	72.22	92.86	86.67	92.31	72.22	81.25	85.71	100.00
	recall	82.35	88.24	76.47	82.35	76.47	76.47	58.82	76.47	76.47	76.47	70.59	76.47	76.47	70.59	82.35
	f1 score	82.35	88.24	72.22	77.78	78.79	81.25	68.97	74.29	83.87	81.25	80.00	74.29	78.79	77.42	90.32
7	precision	76.47	65.00	92.86	72.22	85.71	76.47	80.00	68.75	77.78	73.33	68.42	68.42	68.42	86.67	100.00
	recall	76.47	76.47	76.47	76.47	70.59	76.47	70.59	64.71	82.35	64.71	76.47	76.47	76.47	76.47	76.47
	f1 score	76.47	70.27	83.87	74.29	77.42	76.47	75.00	66.67	80.00	68.75	72.22	72.22	72.22	81.25	86.67
8	precision	92.86	81.25	80.00	93.75	92.86	80.00	66.67	76.47	66.67	82.35	86.67	81.25	76.47	65.00	81.25
	recall	76.47	76.47	70.59	88.24	76.47	70.59	82.35	76.47	82.35	82.35	76.47	76.47	76.47	76.47	76.47
	f1 score	83.87	78.79	75.00	90.91	83.87	75.00	73.68	76.47	73.68	82.35	81.25	78.79	76.47	70.27	78.79

FIGURE 7 Performances of best models for future work citations.

## 4.2.2 | Main path network extraction

The semantic citation networks we analyzed have many small strongly connected components (SCC), so we applied the Simple Search Path Count approach (Jiang et al., 2020), an extension of SPC to deal with cyclic citation networks, for MPN extraction. Their JMPA package<sup>6</sup> (Java package for MPA) was used for implementation. Following Jiang et al. (2020), we segmented the network under analysis to several time slices, extracted top- $K$  ( $K = 10$ ) key-route main paths (Liu et al., 2014) from each slice, and merged them into an MPN. More details are given in Supplementary Section B.2.

## 5 | QUALITATIVE ANALYSIS

For experimental analysis, citation data came from the 2015 version of ACL anthology network (AAN; Radev et al., 2013) about computational linguistics/natural language. Three areas were selected: natural language parsing<sup>7</sup> (AANPar), automatic document summarization (AANSum), and machine translation (AANMT). Due to space limit, this section showcases on AANPar and AANSum to demonstrate the superiority of semantic MPA. Table S4 summarizes the statistics of the (semantic) citation networks and their time slices. The experimental setup is detailed in Supplementary Section B. Key-route MPA was used for main path extraction. It was valid to follow the common practice in MPA to extract semantic MPNs from the largest connected component (CC). This is because all other CCs are all small islands smaller than 2 in the citation network: 91, 191, and 207 in AANSum, AANPar, and AANMT respectively.

## 5.1 | Case Study 1: Natural language parsing

### 5.1.1 | Main path network

For comparison purpose, Figure 8 presents the MPN extracted from the original citation network AANPar. Topic branches are numbered. Seminal papers (verified according to the authors' knowledge about the domain) are in red rectangles, while survey-style papers are in ovals, such as special issue or shared task introduction papers. Table 4 shows a subset of representative main path papers on each topic branch and Table S5 presents the complete list. Topic keywords and short excerpts for certain papers are to assist understanding. Branch 1 describes the early studies about various grammatical formalisms,<sup>8</sup> such as *categorical grammar*, *unification grammar*, *categorical unification grammar*, and *Lambek calculus*. However, since late 1980s, the domain started to have a sense of probabilistic thinking (Branch 2). Branch 3 shows the important development where *Penn Tree-Bank* (J93-2004 and H94-1020) appeared as the most important linguistic resource that most future papers used for developing and evaluating parsing techniques.

Branch 4 represents the mainstream of statistical parsing in the 1990s and 2000s, such as *maximum entropy modeling* (W96-0213, A00-2018) or in another name *log-linear model* (P04-1014), *conditional random fields* (N03-1028), and *max-margin parsing* (W04-3201, P05-1012). Note that, C00-1011 and P00-1009 were two papers on *data-oriented parsing (DOP)* promoted by Rens Bod, which however ceased in the wave of statistical parsing dominated by other proposals presented above. Early studies about *dependency analysis* blossomed into the huge Branch 5 and became the dominant trend since around 2005, further expediated by two important shared tasks W06-2920 and

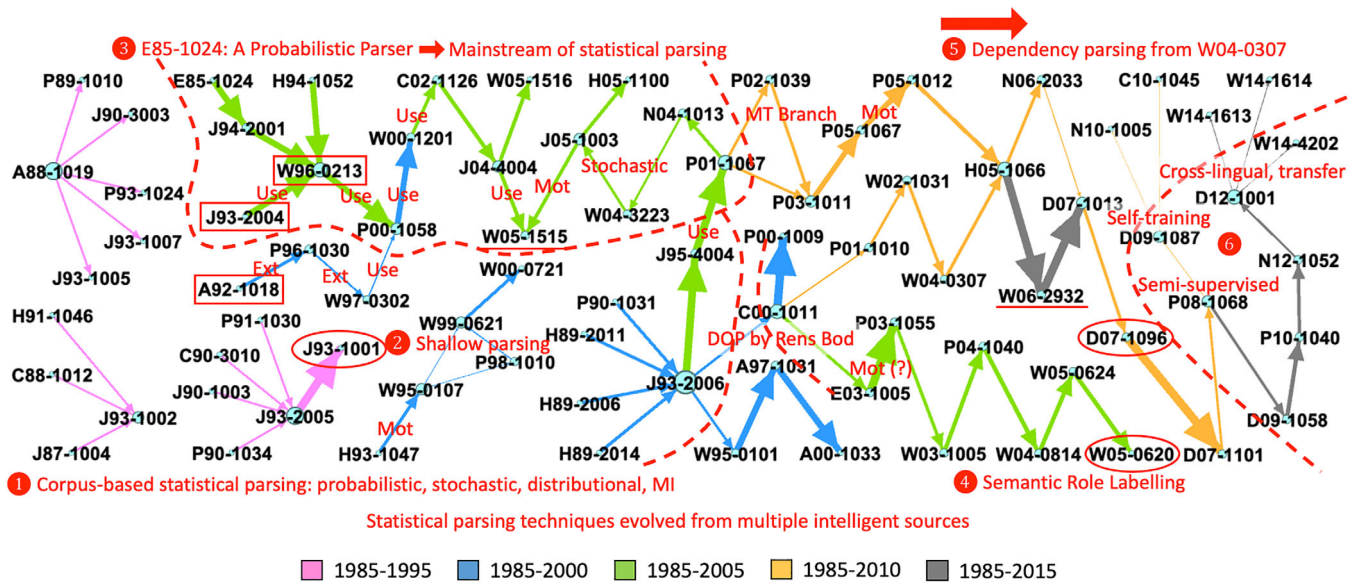


FIGURE 8 Main path network extracted from AANPar.

D07-1096, which then diverted into Branch 6 about dependency parsing of *morphologically rich languages* and Branch 7 about *cross-lingual dependency parsing*. An issue was that many main path papers were connected by incidental citations. For instance, the citation from A00-2018 said that C00-1011 “stays behind the scores of” the former, a weak citation about performance comparison. For another instance, H91-1037 received only 10 citations in our dataset. SPC (H91-1037, J93-2004) was high only because of high-impact citing paper J93-2004 (1,006 citations), although the citation was incidental.

### 5.1.2 | Semantic main path network: Add extension and motivation citations

The above observations motivated us to exploit the semantic relationships between papers in MPA. Figures 9–12 show the semantic MPNs extracted from the four semantic citation networks induced from AANPar, namely AANPar\_add\_Ext\_Mot, AANPar\_plus\_add\_Use, AANPar\_plus\_add\_Sim, and AANPar\_del\_Bkg\_Fut. Interesting chemical reactions occurred when MPA met citation function classification. Each semantic MPN revealed some novel branches or new papers. They collectively drew a more comprehensive picture of domain development. Supplementary Section D presents selected citation context excerpts to help readers understand the citation functions marked on certain edges.

On AANPar\_add\_Ext\_Mot (Figure 9 and Tables 5 and S6 for a complete list of main path papers), the early development of parsing technology was tested. Branch 2 is a new branch about old parsers such as *shift-reduce*

*parsing*, *left-corner parsing*, *tabular parsing*, and *left-to-right (LR) parsing* and so on. Similarly, we saw another (isolated) early development of *probabilistic* approaches (Branch 3; details in Table S6). In addition to A00-2018 as the source of the statistical parsing mainstream, a third source started from E85-1024 (“A probabilistic parser”) to J94-2001 (“Tagging English Text with a Probabilistic Model”) and W96-0213, then through P02-1034 into the new Branch 4 about *multiple parse ranking and re-ranking*. Note that Branch 5 started went into a “dead” end about “Chinese TreeBank” (W00-1201).

From the right part of Figure 9, we saw a branch of *DOP* papers published by Rens Bod until P01-1010. Similar to the evolution pathway in Figure 8, it was gradually merged into the dominant dependency parsing branch. D08-1059 (“A Tale of Two Parsers: Investigating and Combining Graph-based and Transition-based Dependency Parsing”) was motivated (denoted by “Mot” on the edge) by two papers P07-1050 (“K-best Spanning Tree Parsing”) and D07-1013 (“Characterizing the Errors of Data-Driven Dependency Parsing Models”).

Note that, there was a potentially problematic Branch 8 about machine translation (MT) using dependency parsing. Concerning (P05-1012, H05-1066), the citation context excerpt below reveals that although “improving upon” may indicate an extension, the whole context may be recognized as “Similar” or “CoCoGM.” This shows that multilabel classification might be a promising future direction to explore (Lauscher et al., 2022).

“We mentioned above that our approach appears to be similar to that of reranking for statistical parsing (Collins, 2000; Charniak

TABLE 4 Representative main path papers extracted from AANPar.

ACLID	Title
Branch 1	
C86-1045	<i>Categorial Unification Grammars</i>
P87-1012	A Lazy Way To Chart-Parse With <i>Categorial Grammars</i>
C90-2030	Normal Form Theorem Proving For The <i>Lambek Calculus</i>
Branch 2	
H92-1026	Towards History-Based Grammars: Using Richer Models For <i>Probabilistic Parsing</i>
E93-1040	Parsing The Wall Street Journal With The Inside-Outside Algorithm Excerpt: We report grammar inference experiments on partially parsed sentences taken from the <i>Wall Street Journal corpus</i> using the inside-outside algorithm for <i>stochastic context-free grammars</i> .
Branch 3	
J93-2004	<i>Building A Large Annotated Corpus Of English: The Penn Treebank</i>
H94-1020	<i>The Penn Treebank: Annotating Predicate Argument Structure</i>
Branch 4	
A00-2018	A <i>Maximum-Entropy-Inspired</i> Parser
C00-1011	Parsing With The Shortest Derivation (about <i>DOP</i> by Rens Bod)
P00-1009	An Improved Parser For <i>Data-Oriented</i> Lexical-Functional Analysis (about <i>DOP</i> by Rens Bod)
N03-1028	Shallow Parsing With <i>Conditional Random Fields</i>
P04-1014	Parsing The WSJ Using <i>CCG</i> And <i>Log-Linear Models</i>
W04-3201	<i>Max-Margin Parsing</i>
Branch 5	
W06-2920	<i>CoNLL-X Shared Task On Multilingual Dependency Parsing</i>
D07-1096	<i>The CoNLL 2007 Shared Task on Dependency Parsing</i>
D07-1014	Probabilistic Models of <i>Nonprojective Dependency Trees</i>
Branch 6	
W10-1401	Statistical Parsing of <i>Morphologically Rich</i> Languages (SPMRL) What How and Whither
W10-1410	Lemmatization and Lexicalized Statistical Parsing of <i>Morphologically-Rich</i> Languages: the Case of <i>French</i>
Branch 7	
N12-1052	<i>Cross-lingual</i> Word Clusters for Direct <i>Transfer</i> of Linguistic Structure
N13-1126	<i>Target Language Adaptation</i> of Discriminative <i>Transfer</i> Parsers

and Johnson, 2005). While it is true that we are improving upon the output of the automatic parser, we are not considering multiple alternate parses.”

Vague cases exist, such as (W00-1201, C02-1126), a self-citation by D. M. Bikel and D. Chiang. From the citation context excerpt below, expressions like “starting from” and “we have modified” might have been selected as strong signals for extension class (“Ext”).

“The third experiment was on the Chinese Treebank, starting with the same head rules used in (Bikel and Chiang, 2000). These rules were originally ..., and although we have modified them for parsing, ...”

### 5.1.3 | Semantic main path network: Further add usage and similarity citations

By further adding usage citations, that is, on AANPar\_plus\_add\_Use, we saw drastically richer diversity in the development branches (Figure 10, Tables 6 and S7). Again, statistical parsing techniques evolved from multiple intelligent sources (Branches 1-3). A clear notion of “*corpus-based*” parsing emerged (Branch 1). Branch 2 was motivated by H93-1047 (“Automatic Grammar Induction And Parsing Free Text: A Transformation-Based Approach,” a duplicate of P93-1035) and developed into “*shallow parsing*” of words into “*text chunks*.”<sup>9</sup> This time, the seminal paper J93-2004 about the *Penn Treebank* project emerged in Branch 3 and developed through W96-0213 to J04-4004. Most subsequent papers used *Peen Treebank* for development and evaluation. We also saw





TABLE 5 Representative main path papers extracted from AANPar\_add\_Ext\_Mot.

ACLID	Title
Branch 2	
P91-1014	Polynomial Time And Space <i>Shift-Reduce Parsing Of Arbitrary Context-Free Grammars</i>
E93-1036	Generalized <i>Left-Corner Parsing</i>
P94-1017	An Optimal <i>Tabular Parsing Algorithm</i>
Branch 4	
W97-0302	Global Thresholding and <i>Multiple-Pass Parsing</i>
P02-1034	New <i>Ranking Algorithms</i> For Parsing And Tagging: Kernels Over Discrete Structures And The Voted Perceptron
P05-1022	Coarse-To-Fine <i>N-Best Parsing</i> And MaxEnt <i>Discriminative Reranking</i>
...	
Branch 5	
C92-2065	<i>Probabilistic Tree-Adjoining Grammar</i> As A Framework For <i>Statistical Natural Language Processing</i>
C92-2066	<i>Stochastic Lexicalized Tree-Adjoining Grammars</i>
...	
W00-1201	Two Statistical Parsing Models Applied To The <i>Chinese Treebank</i>
Branch 6 (a dead branch)	
C92-3126	A Computational Model Of Language Performance: <i>Data Oriented Parsing</i>
P97-1021	A <i>DOP Model</i> For Semantic Interpretation
Branch 8 (A “wrong” branch)	
P01-1067	A <i>Syntax-Based Statistical Translation Model</i>
P03-1011	Loosely <i>Tree-Based Alignment</i> For <i>Machine Translation</i>
P05-1067	<i>Machine Translation</i> Using Probabilistic Synchronous <i>Dependency Insertion Grammars</i>

which was heavily cited (387 times). The following citation context excerpt proved that similarity citation is indeed relevant to knowledge flow of scientific ideas.

“The maximum entropy models used here are similar in form to those in (Ratnaparkhi, 1996; Berger, Della Pietra, and Della Pietra, 1996; Lau, Rosenfeld, and Roukos, 1993).”

The domain then evolved to the dominant dependency parsing branch (Branch 3), where we were excited to see two new shared tasks about *joint syntactic and semantic dependency parsing* (W08-2121, W09-1201), and then to Branch 4 of subsequent

TABLE 6 Representative Main Path Papers Extracted from AANPar\_plus\_add\_Use.

ACLID	Title
Branch 1	
J93-1002	Generalized <i>Probabilistic LR Parsing</i> Of Natural Language ( <i>Corpora</i> ) With Unification-Based Grammars
J93-1001	<i>Introduction To The Special Issue On Computational Linguistics Using Large Corpora</i>
P90-1031	<i>Parsing The LOB Corpus</i>
Branch 2	
W95-0107	<i>Text Chunking</i> using Transformation-Based Learning
W00-0721	<i>Shallow Parsing</i> by Inferencing with Classifiers
Branch 3	
...	
J93-2004	<i>Building A Large Annotated Corpus Of English: The Penn Treebank</i>
...	
W00-1201	Two <i>Statistical Parsing Models</i> Applied To The <i>Chinese Treebank</i>
...	
P01-1067	A <i>Syntax-Based Statistical Translation Model</i>
Branch 4	
W05-0620	<i>Introduction To The CoNLL-2005 Shared Task: Semantic Role Labeling</i>
Branch 6 (Extended branch about cross-lingual dependency parsing)	
P08-1068	Simple <i>Semi-supervised</i> Dependency Parsing
D09-1087	<i>Self-Training PCFG Grammars</i> with Latent Annotations Across Languages
W14-1613	Distributed Word Representation Learning for <i>Cross-Lingual Dependency Parsing</i>

studies on *semantic dependency parsing* (W09-1208, D09-1004).

#### 5.1.4 | Semantic main path network: Delete neutral and future work citations

Finally, on AANPar\_del\_Bkg\_Fut (Figure 12, Tables 8 and S9), we observed some interesting branches or papers. Since P08-1068, the domain diverted into a new branch about *optimization techniques* used in parsing algorithms, such as *dynamic programming*, *integer linear programming* and *dual decomposition* (Branch 2). Branch 3 was a similar cross-lingual dependency parsing branch, but it evolved into Branch 4 about parsing morphologically rich languages through a new *shared task*

**TABLE 7** Representative main path papers extracted from AANPar\_plus\_add\_Sim.

ACLID	Title
Branch 1	
P99-1047	A Decision-Based Approach To <i>Rhetorical Parsing</i>
J00-3005	The <i>Rhetorical Parsing</i> Of Unrestricted Texts: A Surface-Based Approach
Branch 2	
P02-1042	Generative Models For <i>Statistical Parsing</i> With <i>Combinatory Categorical Grammar</i>
P04-1014	Parsing The WSJ Using <i>CCG</i> And <i>Log-Linear Models</i>
C04-1180	Wide-Coverage Semantic Representations From A <i>CCG Parser</i>
Branch 3 (extended branch of dependency parsing)	
...	
W08-2121	The <i>CoNLL 2008 Shared Task on Joint Parsing of Syntactic and Semantic Dependencies</i>
W09-1201	The <i>CoNLL-2009 Shared Task: Syntactic and Semantic Dependencies in Multiple Languages</i>
Branch 4 (extended to semantic dependency parsing)	
W09-1208	<i>Multilingual Dependency Learning: A Huge Feature Engineering Method to Semantic Dependency Parsing</i>
D09-1004	<i>Semantic Dependency Parsing of NomBank and PropBank: An Efficient Integrated Approach via a Large-scale Feature Selection</i>

(W13-4917), thus provided a complementary view to Branch 6 in Figure 8. We postulate the result is meaningful since dependency parsing was directed by important shared tasks. Note that, deleting neutral and future work citations might result in weaker semantic coherence than by adding more significant citations like extension and similarity (quantified in Section 6.3). For example, N07-1069 only made a result comparison with W06-2928, therefore it is less confident to say scientific ideas flew through this path.

“Here we can compare directly with the best systems for this dataset in CoNLL-X. The best system (Corston-Oliver & Aue, 2006), ....”

In summary, we conjecture that multiple semantic MPNs extracted from different types of semantic citation networks reveal complimentary views and novel knowledge flows, thus should be merged into a more comprehensive representation of scientific domain's topic evolution.

**TABLE 8** Representative main path papers Extracted from AANPar\_del\_Bkg\_Fut.

ACLID	Title
Branch 2	
W08-2102	TAG, <i>Dynamic Programming</i> , and the Perceptron for Efficient, Feature-Rich Parsing
P09-1039	Concise <i>Integer Linear Programming</i> Formulations for Dependency Parsing
D10-1001	On <i>Dual Decomposition</i> and <i>Linear Programming</i> Relaxations for Natural Language Processing
Branch 3	
W06-2928	Dependency Parsing With Reference To Slovene <i>Spanish</i> And <i>Swedish</i>
D07-1119	<i>Multilingual</i> Dependency Parsing and <i>Domain Adaptation</i> using DeSR
P13-2017	Universal Dependency Annotation for <i>Multilingual</i> Parsing
Branch 4	
W13-4917	<i>Overview of the SPMRL 2013 Shared Task: A Cross-Framework Evaluation of Parsing Morphologically Rich Languages</i>
W13-4905, W13-4906 and W13-4910 are all SPMRL 2013 Shared Task papers	

## 5.2 | Case Study 2: Automatic document summarization

Due to space limit, an informative summary is presented here (Figure 13–17). See Tables S10–S14 in Supplementary Section E for the details of main path papers and Supplementary Section F for citation context excerpts. The MPN extracted from AANSum (Figure 13) covered a few early summarization studies centering around the usage of *semantic coherence* devices (Branch 1), such as *discourse structure*, *rhetorical relations*, and *lexical chains* (W97-0703: Using Lexical Chains For Text Summarization), and so on. Then the main body of literature focused on *multidocument summarization* (Branch 2) pioneered by the seminal journal article J98-3005 (“Generating Natural Language Summaries From Multiple On-Line Sources”). The subsequent studies in this topic eventually gave birth to an important *Special Issue* on Summarization (J02-4001). Since the advent of PageRank in 1998, the *graph-based ranking* idea was introduced to the summarization domain for sentence ranking for extractive summarization (Branch 3). Seminal works included P04-3020 (“Graph-Based Ranking Algorithms For Sentence Extraction Applied To Text Summarization”), W04-3252 (“TextRank: Bringing Order Into Texts”),



FIGURE 15 Main path network extracted from AANSum\_plus\_add\_Use.

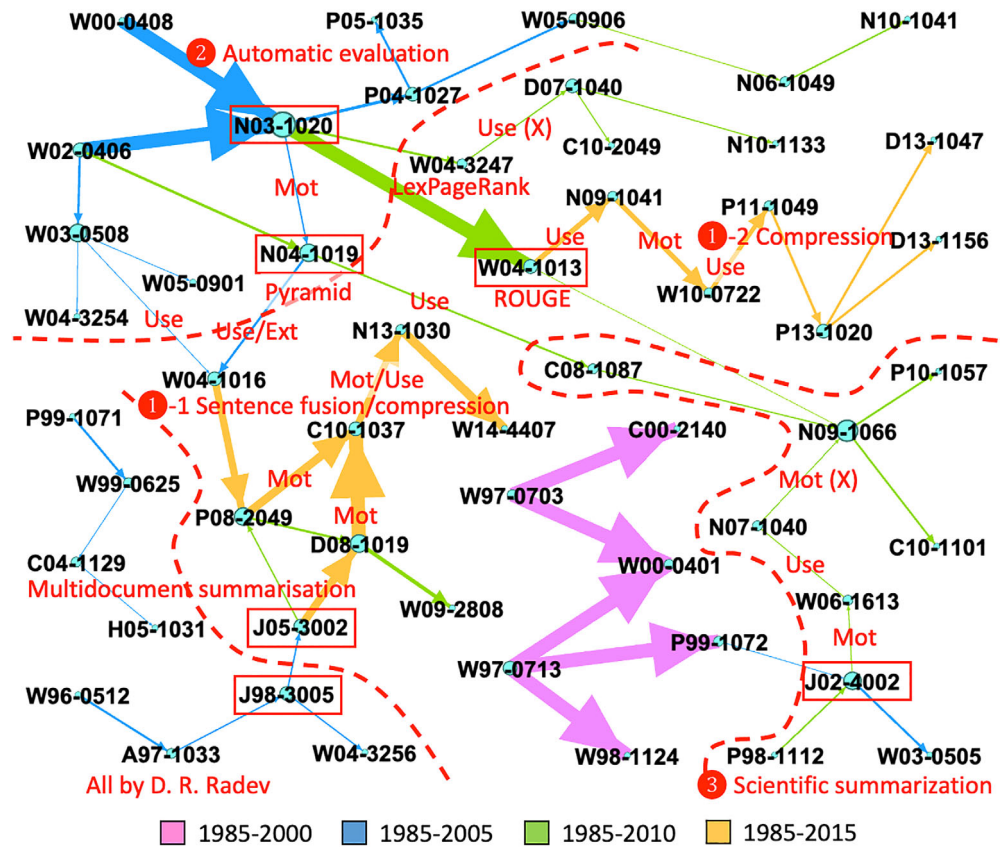
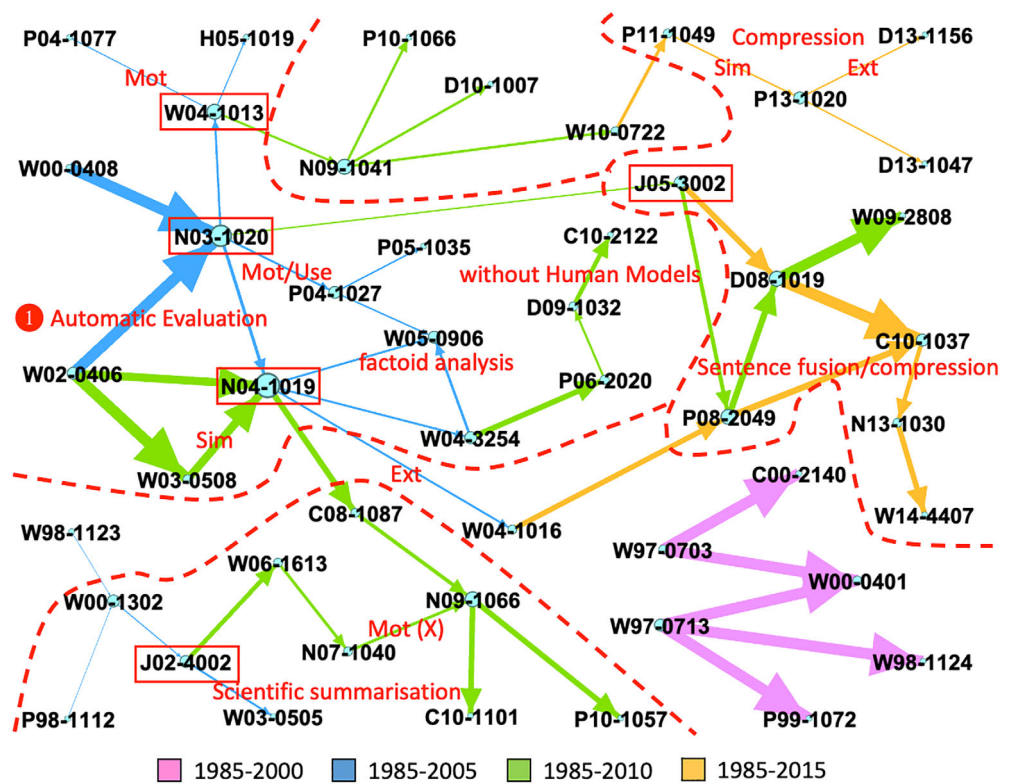


FIGURE 16 Main path network extracted from AANSum\_plus\_add\_Sim.



“Our methodology builds and extends the Teufel and Moens (Teufel and Moens, 2002) approach to automatic summarization.”

In addition to the common topics like multidocument summarization (Branch 2) and graph-based ranking algorithms (Branch 5), we were also excited to see

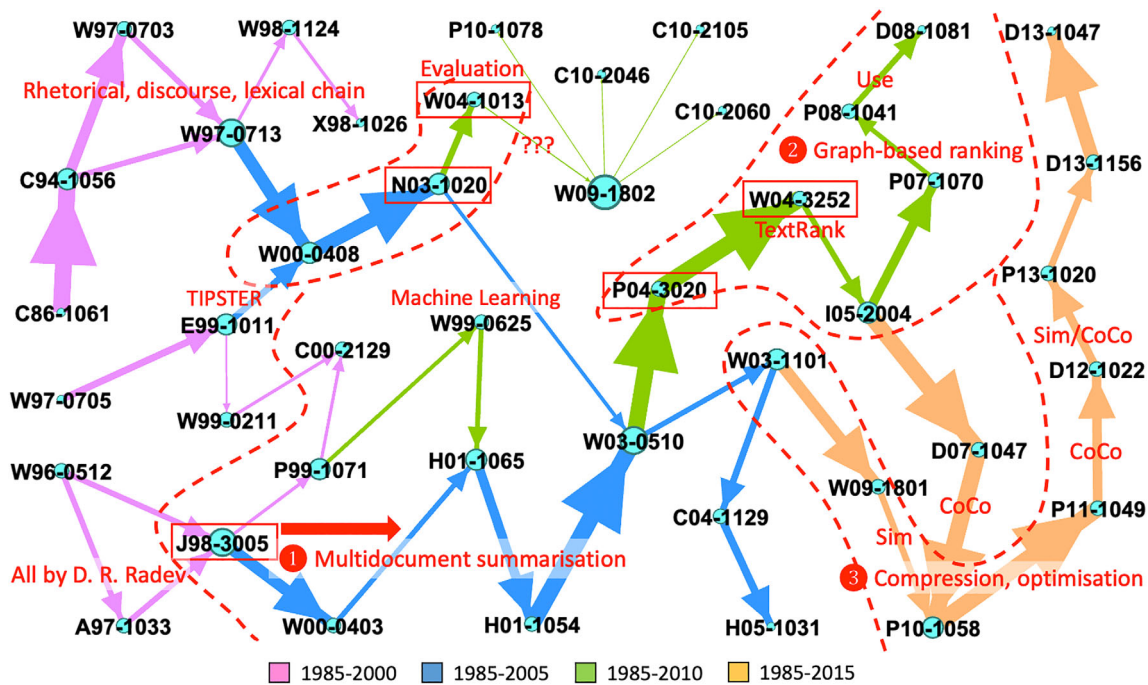


FIGURE 17 Main path network extracted from AANSum\_del\_Bkg\_Fut.

Branch 3 about *automatic evaluation* and related studies. Heavily cited ones included N03-1020 and W04-1013 about the ROUGE package. We also saw more studies about *sentence reduction, compression and fusion* for summarization. Both Branch 4-1 and 4-2 were pioneered by K. R. McKewon in A00-1043 (“Sentence Reduction For Automatic Text Summarization”), A00-2024 (“Cut and Paste Based Text Summarization”), and J05-3002 (“Sentence Fusion For Multidocument News Summarization”).

By further adding usage citations (Figure 15), although we lost the graph-based ranking branch (despite that we got a new paper W04-3247 about LexPageRank), we could uncover more novel topics and branches. Branch 2 about automatic evaluation included more important papers such as N04-1019 about the *Pyramid method* (“Evaluating Content Selection In Summarization: The Pyramid Method”). A significant new branch was Branch 3 about *scientific summarization* at right bottom, starting from the seminal paper J02-4002 to *citation function classification* (W06-1613, N07-1040) and *citation-based summarization* (C08-1087, N09-1066, P10-1057, and C10-1101). By further adding similarity citations (Figure 16), we could see one obvious expansion of Branch 1 about evaluation, starting from *factoid analysis* (W04-3254) to summarization *evaluation without human models*, including D09-1032 (“Automatically Evaluating Content Selection in Summarization without Human Models”) and C10-2022 (“Multilingual Summarization Evaluation without Human Models”), both written by

famous researchers in this domain (A. Nenkova and H. Saggion respectively).

Finally, the MPN extracted from AANSum\_del\_Bkg\_Fut (Figure 17) recovered the vanished or shrunk branches about *multidocument summarization* (Branch 1) and *graph-based ranking* (Branch 2), and at the same time introduced some new papers, such as C04-1129 for Branch 1 (“Syntactic Simplification For Improving Content Selection In Multi-Document Summarization”), P08-1048 for Branch 2 (“Summarizing Emails with Conversational Cohesion and Subjectivity,” whose abstract says “Second, we use two graph-based summarization approaches, ..., to extract sentences as summaries.”), and W09-1802 (“A Scalable Global Model for Summarization,” whose abstract says “We present an Integer Linear Program for ... for automatic summarization.”) and C10-2105 (“Opinion Summarization with Integer Linear Programming Formulation for Sentence Extraction and Ordering”) for Branch 3 about optimization methods for summarization.

Again, by gradually adding more citation semantics, the semantic MPNs together proved to be more expressive than the semantics-agnostic counterpart.

## 6 | QUANTITATIVE ANALYSIS

Few studies touched quantitative MPA evaluation. Filipin (2021) claimed that it is questionable if a main path is representative of the real technological trajectory because,

based on domain experts' opinions, main path may be "limited to a much narrower neighborhood of the technology space than it really is" and may miss many crucial studies and big players of the analyzed field. Huang et al. (2022) claimed to have achieved better convergence, which was only qualitatively justified. The current situation called us to propose a three-way framework for quantitative MPA evaluation. The first drawback pointed out by Filippin implies that a good main path should have a good *coverage* of the scientific topics of an analyzed domain. It should also include as many critical studies as possible. We name this aspect the *pertinence* of main path. Furthermore, according to Huang et al., nearby main path nodes should exhibit a certain level of local clustering and show higher topical *coherence*. Our framework evaluated all these three aspects.

## 6.1 | Topic modeling

Coverage and coherence were both defined based on topic modeling, here LDA (Blei et al., 2003) trained using the Gensim package.<sup>11</sup> Each article  $u$  in the citation network, denoted as  $CN$ , was represented by its topic distribution  $\mathbf{u} = [u_1, \dots, u_t, \dots, u_T]$ , where  $T$  is topic number,  $u_t$  is the probability of article  $u$  belonging to topic  $t$ , and  $\sum_{t=1}^T u_t = 1$ . Two issues arose: the right value of  $T$  and the right number of training epochs  $P$  (to avoid overfitting LDA training). Supplementary Section G details how to decide these values. In summary, we trained several LDA models with a range of values of  $T$  for evaluation and reported the average. For AANPar,  $T$  values fell in {10, 11, ..., 20, 22, 24, 26}. For AANSum, and AANMT, the maximum value of  $T$  was set to 20. The right value of  $P$  was set to 50, 40, and 50 for AANPar, AANSum, and AANMT respectively.

## 6.2 | Topical coverage

Let  $MN$  denote an extracted MPN. *Topical coverage* measures how well  $MN$  covers the topics of the analyzed domain. It is approximated by the closeness between the topic distribution of  $MN$ , denoted as  $dist_{tpk}(MN)$ , and the topic distribution of  $CN$ , denoted as  $dist_{tpk}(CN)$ , both of which are averaged over the enclosed publications. In evaluation, we used Hellinger distance to measure topical coverage, defined below:

$$cov_{tpk}(MN, CN) = D_{Hellinger}(dist_{tpk}(MN), dist_{tpk}(CN)), \quad (1)$$

where the Hellinger distance between two vectors  $\mathbf{u}$  and  $\mathbf{v}$  is defined as

$$D_{Hellinger}(\mathbf{u}, \mathbf{v}) = \frac{1}{\sqrt{2}} \sqrt{\sum_i (\sqrt{u_i} - \sqrt{v_i})^2}. \quad (2)$$

The smaller the Hellinger distance is, the better topical coverage is in our sense. Table 9 shows the results. Each "Δ%" column shows the difference of the corresponding semantic MPN from the vanilla MPN in percentage format. Thus, a positive percentage means a decrease in topical coverage and a negative percentage means increase. The upward and downward arrows signify a further increase and decrease from the semantic MPN in the column to the left. On all three datasets, compared to the semantics agnostic counterpart (the "MPN" column), topical coverage decreased (signified by upward arrows) by adding extension and motivation citations (the "add\_Ext\_Mot" column), but adding usage relations lead to improved topical coverage (signified by downward arrows in the "plus\_add\_Use" column). This is meaningful because publications linked with extension and motivation citations are technically closer. On the contrary, usage can be about a variety of different things, from algorithm and method to data and definition, and so on, and thus results in main paths that are topically more diverse. Two composite semantic MPNs were extracted: "add\_Combined" corresponds to the composite semantic MPN which merged three semantic MPNs corresponding to "add\_Ext\_Mot," "plus\_add\_Use" and "plus\_add\_Sim"; "del\_Combined" corresponds to the composite semantic MPN which further merged the semantic MPN corresponding to "del\_Bkg\_Fut." The results proved that different types of semantic MPNs complemented each other and collectively worked better, that is, covering and approximating the topic distribution of the underlying domain much better. Meanwhile, we also confess that better coverage was partially because composite semantic MPNs were larger in size (also see Table 11).

## 6.3 | Topical coherence

A perfect definition of coherence does not exist. We tried to analyze coherence by adapting the coherence definition originally proposed to evaluate topic model quality (Newman et al., 2010, p. 102). Given a main path network  $MN$ , we defined *topical coherence* as the mean of distances between all pairs of main path nodes:

$$coh_{tpk}(MN) = mean\{D(u, v), \forall (u, v) \in MN\}, \quad (3)$$

where  $D(u, v)$  is the distance between the topic distributions of  $u$  and  $v$ . Again, Hellinger distance defined in Eq. (2) was used.

**TABLE 9** Topical coverage of main path networks.

	MPN	add_Ext_Mot		plus_add_Use		plus_add_Sim		add_Combined		del_Bkg_Fut		del_Combined	
		<i>cov<sub>tpk</sub></i>	<i>cov<sub>tpk</sub></i>	$\Delta\%$	<i>cov<sub>tpk</sub></i>	$\Delta\%$	<i>cov<sub>tpk</sub></i>	$\Delta\%$	<i>cov<sub>tpk</sub></i>	$\Delta\%$	<i>cov<sub>tpk</sub></i>	$\Delta\%$	<i>cov<sub>tpk</sub></i>
AANSum	0.0611	0.0647	+6.79% ↑	0.0591	-0.87% ↓	0.0679	+13.92% ↑	0.0509	-15.20% ↓	0.0630	+5.25%	0.0441	-26.53% ↓
AANPar	0.0582	0.0700	+25.13% ↑	0.0496	-8.40% ↓	0.0420	-24.34% ↓	0.0387	-29.62% ↓	0.0694	+21.57%	0.0380	-32.43% ↓
AANMT	0.0696	0.0794	+24.78% ↑	0.0617	-2.34% ↓	0.0697	+9.34% ↑	0.0621	-2.08% ↓	0.0619	-3.93%	0.0497	-20.18% ↓

**TABLE 10** Topical coherence of main path networks.

	MPN	add_Ext_Mot		plus_add_Use		plus_add_Sim		add_Combined		del_Bkg_Fut		del_Combined		
		<i>coh<sub>tpk</sub></i>	<i>coh<sub>tpk</sub></i>	$\Delta\%$	<i>coh<sub>tpk</sub></i>	$\Delta\%$	<i>coh<sub>tpk</sub></i>	$\Delta\%$	<i>coh<sub>tpk</sub></i>	$\Delta\%$	<i>coh<sub>tpk</sub></i>	$\Delta\%$	<i>coh<sub>tpk</sub></i>	$\Delta\%$
Evaluate on <i>MN</i>	AANSum	0.5518	0.5350	-3.18%	0.5456	-1.30% ↑	0.5428	-1.70% ↓	0.5423	-1.84% ↓	0.5505	-0.26%	0.5484	-0.67% ↓
	AANPar	0.4504	0.4448	-1.34%	0.4600	+2.14% ↑	0.4504	-0.05% ↓	0.4488	-0.40% ↓	0.4472	-0.71%	0.4484	-0.48% ↑
	AANMT	0.4327	0.4261	-1.41%	0.4394	+1.61% ↑	0.4138	-4.43% ↓	0.4246	-1.77% ↑	0.4299	-0.70%	0.4266	-1.38% ↓
Evaluate on <i>CN[MN]</i>	AANSum	0.5709	0.5736	+0.51%	0.5642	-1.25% ↓	0.5529	-3.17% ↓	0.5631	-1.39% ↑	0.5720	+0.31%	0.5698	-0.16% ↓
	AANPar	0.4748	0.4602	-3.00%	0.4878	+2.79% ↑	0.4791	+0.95% ↓	0.4730	-0.31% ↓	0.4718	-0.61%	0.4726	-0.40% ↑
	AANMT	0.4492	0.4529	+0.85%	0.4576	+1.96% ↑	0.4489	-0.02% ↓	0.4535	+1.02% ↑	0.4461	-0.70%	0.4545	+1.20% ↑



“add\_Combined” and “del\_Combined” rows against the “MPN” row, the recalls of the former were more than doubled on AANPar and AANSum, and gained more than 65% relative increase on AANMT. Recall that, it is extremely important that as many crucial studies as possible are detected by MPA. At the same time, F1 scores were also largely improved except on AANMT\_add\_Combined. In addition, from the last three rows, we saw that “add\_Combined” and “del\_Bkg\_Fut” results also complemented each other. The most extreme case was on AANMT: the sum of recalls of “add\_Combined” and “del\_Bkg\_Fut” was only slightly larger than the recall of “del\_Combined,” implying that they returned drastically different subsets of gold standards. This justifies our claim that semantic MPNs may exhibit higher diversity to complement each other, and it would be better to merge them for a more comprehensive view. Finally, the recalls and F1 scores on all three datasets corroborate with the findings of Filippin (2021) about MPA’s unsatisfactory recognition rate of the most significant studies. Although semantic MPA proved to improve ranking pertinence by a large margin, there seemed to still large space to improve recall. To achieve this, we guess that it may be helpful to start and guide main path exploration by first ranking and selecting important publications in some way (Bae et al., 2014; Zhang et al., 2014; Tao et al., 2017; Ding et al., 2022).

## 7 | CONCLUSIONS

This paper advocated a novel semantic main path network analysis approach for extracting the scientific backbone from a citation network based on citation function analysis. First, according to per-class performance analysis, the best models for extension, motivation, usage, similarity, neutral (equiv. background) and future work citations were cherry-picked from 55 contextualized citation function classification models trained from 11 model architectures based on SciBERT. Then, four types of semantic citation networks were created by gradually adding extension and motivation citations, usage citations, and similarity citations in a recall-oriented fashion, and by deleting neutral and future work citations in a precision-oriented way. On each semantic citation network, semantic main path network was extracted by merging the top- $K$  key-route main paths extracted from different time slices of the network. Meanwhile, for the first time, this paper performed quantitative main path analysis evaluation by proposing a three-way framework consisting of topical coverage, topical coherence and ranking pertinence. The effectiveness of semantic main path network analysis was demonstrated on three

computational linguistics fields, namely natural language parsing, automatic text summarization and machine translation.

Qualitative analysis showed that each semantic main path network was able to reveal novel topic branches, new important papers of existing branches, and the development pathways between papers and branches, thus provided complementary views of domain evolution. For example, for large domains such as natural language parsing that were guided by a few seminal studies (like Penn Treebank) and ground-breaking shared tasks, the semantic main path networks were much better at finding these representative works, such as the two early shared tasks on (multilingual) dependency parsing and more future shared tasks on a plethora of topics including semantic dependency parsing, semantic role labeling and dependency parsing of morphologically rich languages, most of which were missed by traditional main path analysis. For automatic text summarization, the semantic main path network approach was able to find an important novel branch about summarization evaluation and the branch about optimization methods for summarization, at the same time enrich the multidocument summarization, graph-based ranking and sentence fusion/compression branches that were recognized by the traditional approach.

Merging multiple semantic main path networks resulted in significantly better topical coverage. When main path analysis is seen as a method to return an unordered set of top-ranked studies, the composite semantic main path networks achieved much better ranking pertinence based on expert-selected gold standards, thus proved to be more comprehensive representations of scientific development. In addition, extension, motivation and similarity citations proved to achieve better semantic coherence on all three datasets than traditional approaches which ignore citation semantics, but adding usage citations may introduce topical diversity, which resulted in lower coherence but higher coverage. In the extracted semantic main path networks, most recognized citation relations were more relevant to uncovering the knowledge flow among scientific ideas. On the contrary, in the traditional approach, many main path papers were connected via incidental citations such as neutral citations. Therefore, we conclude that the semantic main path network analysis approach can discover more pertinent topic branches, uncover more coherent knowledge flows, and provide a more comprehensive scientific domain representation.

## ACKNOWLEDGMENT

Xiaorui Jiang is partially supported by National Office for Philosophy and Social Sciences of China (18ZDA238).

## ORCID

Xiaorui Jiang  <https://orcid.org/0000-0003-4255-5445>

## ENDNOTES

- <sup>1</sup> Refer to Kuan (2023) for more discussions.
- <sup>2</sup> <https://aclanthology.org/>
- <sup>3</sup> Supplementary materials: [https://github.com/xiaoruijiang/CFC\\_MPN/blob/main/jasist2022\\_v2\\_SM\\_for\\_review.pdf](https://github.com/xiaoruijiang/CFC_MPN/blob/main/jasist2022_v2_SM_for_review.pdf)
- <sup>4</sup> This choice was supported by the claim made by Lauscher et al. (2022) that most citation instances' functions could be determined only using citance alone.
- <sup>5</sup> When  $\mathbf{f} = \mathbf{h}$ , depending on context\_type, the number of model variants is 2. When  $\mathbf{f} = [\mathbf{h}; \mathbf{s}]$ , the number of model variants is: 2 (context\_type = "sequential") + 2 × 3 (context\_type = "hierarchical") = 8. When  $\mathbf{f} = [\mathbf{h}; \mathbf{c}]$ , if context\_type = "sequential", the model variant number is 2; otherwise, if context\_type = "hierarchical", it is 3 × 2 = 6 (3 sentence poolers by 2 context encoders). When  $\mathbf{f} = [\mathbf{h}; \mathbf{s}; \mathbf{c}]$ , if context\_type = "sequential", the model variant number is 2 × 2 (2 citance encoders multiplied by 2 context encoders) = 4; otherwise if context\_type = "hierarchical", there are 2 × 3 × 2 = 12 model variants (2 citance encoders by 3 sentence poolers by 2 context encoders). Therefore, there are in total 2 + 8 + (2 + 6) + (4 + 12) = 34 model variants.
- <sup>6</sup> <https://github.com/xiaoruijiang/JMPA>
- <sup>7</sup> Parsing: Parsing, syntax analysis, or syntactic analysis is the process of analyzing a string of symbols, either in natural language, computer languages or data structures, conforming to the rules of a formal grammar. See Wikipedia page: <https://en.wikipedia.org/wiki/Parsing>.
- <sup>8</sup> Note that more grammars were proposed even earlier, outside our time range of analysis.
- <sup>9</sup> From Wikipedia, shallow parsing is also chunking or light parsing: [https://en.wikipedia.org/wiki/Shallow\\_parsing](https://en.wikipedia.org/wiki/Shallow_parsing)
- <sup>10</sup> Both semantic role labeling and dependency parsing became rather standalone topics and had bespoke monographs on these two topics.
- <sup>11</sup> <https://radimrehurek.com/gensim>
- <sup>12</sup> They are available at: [https://github.com/xiaoruijiang/scirank/tree/main/datasets/gold\\_standards/ACL](https://github.com/xiaoruijiang/scirank/tree/main/datasets/gold_standards/ACL). Note that, to construct GS-PAR, we referred to Jiang et al.'s gold standard papers about computational linguistics/natural language (Jiang et al., 2019), and manually picked out the papers about natural language parsing technologies, because the surveys we were able to find could not cover the whole area of natural language processing.

## REFERENCES

- Abu-Jbara, A., Erza, J., & Radev, D. (2013). Purpose and polarity of citation: Towards NLP-based bibliometrics. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT'13)* (pp. 596–606). Association for Computational Linguistics. <https://aclanthology.org/N13-1067>
- Bae, D.-H., Hwang, S.-M., Kim, S.-W., & Faloutsos, C. (2014). On constructing seminal paper genealogy. *IEEE Transactions on Cybernetics*, 44(1), 54–65. <https://doi.org/10.1109/TCYB.2013.2246565>
- Batagelj, V. (2003). *Efficient algorithms for citation network analysis*. <https://arxiv.org/abs/cs/0309023>
- Beltagy, I., Lo, K., & Cohan, A. (2019). SciBERT: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP'19)* (pp. 3615–3620). Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1371>
- Blei, D., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3, 993–1022. <https://www.jmlr.org/papers/v3/blei03a.html>
- Chen, L., Xu, S., Zhu, L., Zhang, J., Xu, H., & Yang, G. (2022). A semantic main path analysis method to identify multiple development trajectories. *Journal of Informetrics*, 12, 101281. <https://doi.org/10.1016/j.joi.2022.101281>
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., & Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41(6), 391–407. [https://doi.org/10.1002/\(SICI\)1097-4571\(199009\)41:6%3C391::AID-ASI%3E3.0.CO;2-9](https://doi.org/10.1002/(SICI)1097-4571(199009)41:6%3C391::AID-ASI%3E3.0.CO;2-9)
- Dietz, L., Cickel, S., & Scheffer, T. (2007). Unsupervised prediction of citation influences. In *Proceedings of the 24th International Conference on Machine Learning (ICML'07)* (pp. 233–240). Association for Computing Machinery. <https://doi.org/10.1145/1273496.1273526>
- Ding, J., Xiang, T., Ou, Z., Zuo, W., Zhao, R., Lin, C., Zheng, Y., & Liu, B. (2022). Tell me how to survey: literature review made simple with automatic reading path generation. In *Proceedings of the 2022 IEEE 38th International Conference on Data Engineering (ICDE'22)* (pp. 3426–3438). IEEE. <https://doi.org/10.1109/ICDE53745.2022.00322>
- Dong, C., & Schäfer, U. (2011). Ensemble-style self-training on citation classification. In *Proceedings of 5th International Joint Conference on Natural Language Processing (IJCNLP'11)* (pp. 623–631). Association for Computational Linguistics. <https://aclanthology.org/I11-1070>
- Filippin, F. (2021). Do main paths reflect technological trajectories? Applying main path analysis to the semiconductor manufacturing industry. *Scientometrics*, 126(8), 6443–6477. <https://doi.org/10.1007/s11192-021-04023-9>
- Garfield, E., Pudovkin, A. I., & Istomin, V. S. (2003). Why do we need algorithmic historiography? *Journal of the American Society for Information Science and Technology*, 54(5), 400–412. <https://doi.org/10.1002/asi.10226>
- Ghosal, T., Tiwary, P., Patton, R., & Stahl, C. (2022). Towards establishing a research lineage via identification of significant citations. *Quantitative Science Studies*, 2(4), 1511–1528. [https://doi.org/10.1162/qss\\_a\\_00170](https://doi.org/10.1162/qss_a_00170)
- Hassan, S.-U., Safder, I., Akram, A., & Kamiran, F. (2018). A novel machine-learning approach to measuring scientific knowledge flows using citation context analysis. *Scientometrics*, 116, 973–996. <https://doi.org/10.1007/s11192-018-2767-x>
- Hernández-Alvarez, M., Gómez, J. M., & Martínez-Barco, P. (2017). Citation function, polarity and influence classification. *Natural Language Engineering*, 23(4), 561–588. <https://doi.org/10.1017/S1351324916000346>
- Huang, C.-H., Liu, J. S., Ho, M. H.-C., & Chou, T.-C. (2022). Towards more convergent main paths: A relevance-based

- approach. *Journal of Informetrics*, 16, 101317. <https://doi.org/10.1016/j.joi.2022.101317>
- Hummon, N. P., & Doreian, P. (1989). Connectivity in a citation network: The development of DNA theory. *Social Networks*, 11(1), 39–63. [https://doi.org/10.1016/0378-8733\(89\)90017-8](https://doi.org/10.1016/0378-8733(89)90017-8)
- Iqbal, S., Hassan, S.-I., Aljohan, N. R., Alelyani, S., Nawaz, R., & Bornmann. (2021). A decade of in-text citation analysis based on natural language processing and machine learning techniques: an overview of empirical studies. *Scientometrics*, 126, 6551–6599. <https://doi.org/10.1007/s11192-021-04055-1>
- Jha, R., Abu-Jbara, A., Qazvinian, V., & Radev, D. R. (2017). NLP-driven citation analysis for scientometrics. *Natural Language Engineering*, 23(1), 93–130. <https://doi.org/10.1017/S1351324915000443>
- Jiang, X., & Zhuge, H. (2019). Forward search path count as an alternative indirect citation impact indicator. *J. Informetrics*, 13(4), 100977. <https://doi.org/10.1016/j.joi.2019.100977>
- Jiang, X., Zhu, X., & Chen, J. (2020). Main path analysis on cyclic citation networks. *Journal of the Association for Information Science and Technology*, 71(5), 578–595. <https://doi.org/10.1002/asi.24258>
- Jurgens, D., Kumar, S., Hoover, R., McFarland, D., & Jurafsky, D. (2018). Measuring the evolution of a scientific field through citation frames. *Transactions of the Association for Computational Linguistics*, 6, 391–406. [https://doi.org/10.1162/tacl\\_a\\_00028](https://doi.org/10.1162/tacl_a_00028)
- Kim, E. H. J., Jeong, Y. K., Kim, Y. H., & Song, M. (2022). Exploring scientific trajectories of a large-scale dataset using topic-integrated path extraction. *Journal of Informetrics*, 16(1), 101242. <https://doi.org/10.1016/j.joi.2021.101242>
- Kim, J., & Shin, J. (2018). Mapping extended technological trajectories: Integration of main paths, derivative paths, and technology junctures. *Scientometrics*, 116(3), 12–17. <https://doi.org/10.1007/s11192-018-2834-3>
- Kim, M., Baek, I., & Song, X. (2018). Topic diffusion analysis of a weighted citation network in biomedical literature. *Journal of the Association for Information Science and Technology*, 69(2), 329–342. <https://doi.org/10.1002/asi.23960>
- Kuan, C.-H. (2020). Regarding weight assignment algorithms of main path analysis and the conversion of arc weights to node weights. *Scientometrics*, 124, 775–782. <https://doi.org/10.1007/s11192-020-03468-8>
- Kuan, C.-H. (2023). Does main path analysis prefer longer paths? *Scientometrics*, 128, 841–851. <https://doi.org/10.1007/s11192-022-04543-y>
- Kunnath, S. N., Herrmannova, D., Pride, D., & Knoth, P. (2022). A meta-analysis of semantic classification of citations. *Quantitative Science Studies*, 2(4), 1170–1215. [https://doi.org/10.1162/qss\\_a\\_00159](https://doi.org/10.1162/qss_a_00159)
- Lauscher, A., Brandon, K., Kuehl, B., Johnson, S., Jurgens, D., Cohan, A., & Lo, K. (2022). MULTICITE: Modeling realistic citations requires moving beyond the single-sentence single-label setting. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL'22)* (pp. 1875–1888). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.naacl-main.137>
- Liu, J. S., Chen, H.-H., Ho, M. H.-C., & Li, Y.-C. (2014). Citations with different levels of relevancy: Tracing the main paths of legal opinions. *Journal of the American Society of Information Science and Technology*, 65(12), 2479–2488. <https://doi.org/10.1002/asi.23135>
- Liu, J. S., & Kuan, C.-H. (2016). A new approach for main path analysis: Decay in knowledge diffusion. *Journal of the Association for Information Science and Technology*, 67(2), 465–476. <https://doi.org/10.1002/asi.23384>
- Liu, J. S., & Lu, L. Y. Y. (2012). An integrated approach for main path analysis: The development of the Hirsch index as an example. *Journal of the American Society for Information Science and Technology*, 59(12), 1948–1962. <https://doi.org/10.1002/asi.21692>
- Liu, J. S., Lu, L. Y. Y., & Ho, M. H.-C. (2019). A few notes on main path analysis. *Scientometrics*, 119, 379–391. <https://doi.org/10.1007/s11192-019-03034-x>
- Liu, J. S., Lu, L. Y. Y., & Ho, M. H.-C. (2020). A note on choosing traversal counts in main path analysis. *Scientometrics*, 124, 783–785. <https://doi.org/10.1007/s11192-020-03469-7>
- Lu, W., Meng, R., & Liu, X. (2014). A deep scientific literature mining-oriented framework for citation content annotation. *Journal of Library Science in China*, 40(214), 93–104. (in Chinese). <https://doi.org/10.13530/j.cnki.jlis.140029>
- Lucio-Arias, D., & Leydesdorff, L. (2008). Main-path analysis and path-dependent transition in HistCite-based historiograms. *Journal of the American Society of Information Science and Technology*, 27(1), 25–45. <https://doi.org/10.1002/asi.20903>
- Lyu, D., Ruan, X., Xie, J., & Cheng, Y. (2021). The classification of citing motivations: a meta-synthesis. *Scientometrics*, 126, 3243–3264. <https://doi.org/10.1007/s11192-021-03908-z>
- Munkhdalai, T., Lalor, J., & Yu, H. (2016). Citation analysis with neural attention models. In *Proceedings of the Seventh International Workshop on Health Text Mining and Information Analysis (LOUHI'16)* (pp. 69–77). Association for Computational Linguistics. <https://doi.org/10.18653/v1/W16-6109>
- Newman, D., Lau, J. H., Grieser, K., & Baldwin, T. (2010). Automatic evaluation of topic coherence. In *Proceedings of Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL'10)* (pp. 100–108). Association for Computational Linguistics. <https://aclanthology.org/N10-1012>
- Radev, D. R., Muthukrishnan, P., Qazvinian, V., & Abu-Jbara, A. (2013). The ACL anthology network corpus. *Language Resource and Evaluation*, 47(4), 919–944. <https://doi.org/10.1007/s10579-012-9211-2>
- Su, X., Prasad, A., Kan, M.-Y., & Sugiyama, K. (2019). Neural multi-task learning for citation function and provenance. In *Proceedings of the 2019 ACM/IEEE Joint Conference on Digital Libraries (JCDL'19)* (pp. 394–395). IEEE. <https://doi.org/10.1109/JCDL.2019.00122>
- Tao, S., Wang, X., Huang, W., Chen, W., Wang, T., & Lei, K. (2017). From citation network to study map: A novel model to reorganize academic literatures. In *In Proceedings of the 26th International Conference on World Wide Web Companion (WWW'17 Companion)* (pp. 1225–1232). ACM. <https://doi.org/10.1145/3041021.3053059>
- Teufel, S. (2010). *The structure of scientific articles: applications to citation indexing and summarization*. Centre for the Study of Language & Information.

- Teufel, S., Siddharthan, A., & Tidhar, D. (2006a). Automatic classification of citation function. In *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing (EMNLP'06)* (pp. 103–110). Association for Computational Linguistics. <https://aclanthology.org/W06-1613>
- Teufel, S., Siddharthan, A., & Tidhar, D. (2006b). An annotation scheme for citation function. In *Proceedings of the 7th SIGdial Workshop on Discourse and Dialogue (SIGdial'06)* (pp. 80–87). Association for Computational Linguistics. <https://aclanthology.org/W06-1312>
- Valenzuela, M., Ha, V., & Etzioni, O. (2015). Identifying meaningful citations. In *Proceedings of the Workshops of Scholarly Big Data: AI Perspectives, Challenges, and Ideas at the 29th AAAI Conference on Artificial Intelligence*. AAAI Press. <https://allenai.org/data/meaningful-citations>
- Verspagen, B. (2007). Mapping technological trajectories as patent citation networks: A study on the history of fuel cell research. *Advances in Complex Systems*, 10(1), 93–115. <https://doi.org/10.1142/S0219525907000945>
- Yu, D., & Pan, T. (2021). Tracing the main path of interdisciplinary research considering citation preference: A case from blockchain domain. *Journal of Informetrics*, 16(2), 101136. <https://doi.org/10.1016/j.joi.2021.101136>
- Yu, D., & Sheng, L. (2021). Influence difference main path analysis: Evidence from DNA and blockchain domain citation networks.

*Journal of Informetrics*, 15(4), 101186. <https://doi.org/10.1016/j.joi.2021.101186>

- Zhang, H., Li, L., Li, T., & Wang, D. (2014). PatentDom: Analyzing patent relationships on multi-view patent graphs. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management (CIKM'14)* (pp. 1369–1378). ACM. <https://doi.org/10.1145/2661829.2662031>

## SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

**How to cite this article:** Jiang, X., & Liu, J. (2023). Extracting the evolutionary backbone of scientific domains: The semantic main path network analysis approach based on citation context analysis. *Journal of the Association for Information Science and Technology*, 74(5), 546–569. <https://doi.org/10.1002/asi.24748>