



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Leibniz
Universität
Hannover



Semantically Orthogonal Framework for Citation Classification: Disentangling Intent and Content

Changxu Duan *Technical University of Darmstadt, Darmstadt, Germany*

Zhiyin Tan *L3S Research Center, Leibniz University Hannover, Hannover, Germany*

The 29th International Conference on Theory
and Practice of Digital Libraries

TPDL
2025

September 23-26, 2025
Tampere, Finland

Background & Motivation

"Citations have long been used to characterize the state of a scientific field and to identify influential works.

However, writers use citations for different purposes, and this varied purpose influences uptake by future scholars."

—— Jurgens et al., *Measuring the Evolution of a Scientific Field through Citation Frames*,
Transactions of the Association for Computational Linguistics, 6, 391-406., 2018

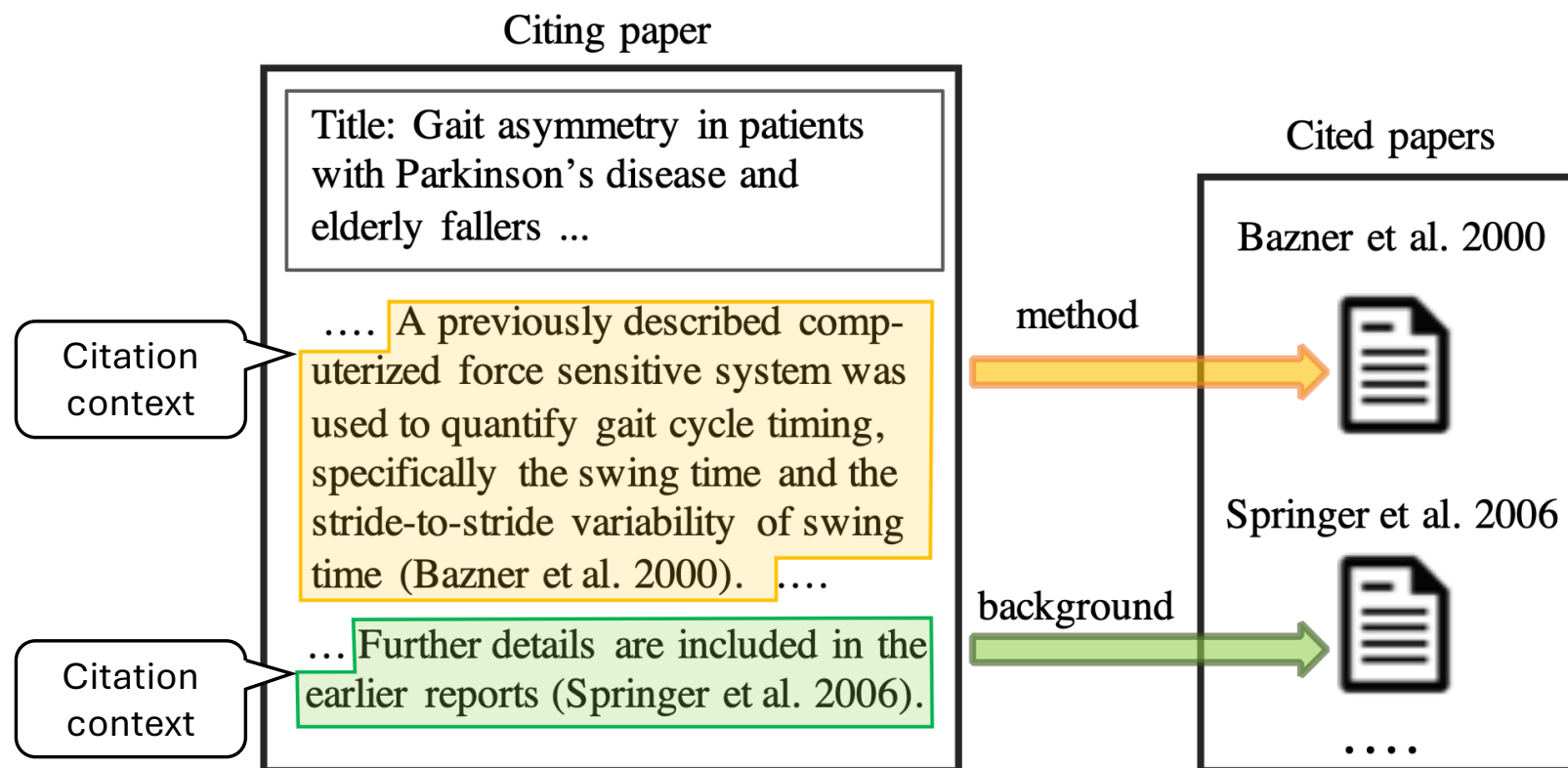
Background & Motivation

*Citations are also typically used as the main **measure** for assessing **impact** of **scientific publications**, **venues**, and **researcher** (Li and Ho, 2008).*

*...identifying the intent of citations is critical in **improving automated analysis of academic literature** and **scientific impact measurement** (Leydesdorff, 1998; Small, 2018).*

—— Cohan et al., *Structural Scaffolds for Citation Intent Classification in Scientific Publications*, *Proceedings of NAACL-HLT*, pages 3586–3596, 2019

Background & Motivation



— Cohan et al., *Structural Scaffolds for Citation Intent Classification in Scientific Publications*, *Proceedings of NAACL-HLT*, pages 3586–3596, 2019

Background & Motivation

The screenshot shows the Semantic Scholar interface for the paper 'GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding'. The page includes a search bar at the top with 229,061,554 papers, a 'Sign In' button, and a 'Create Free Account' button. The paper's DOI is 10.18653/v1/W18-5446 and its Corpus ID is 5034059. The title is prominently displayed, followed by the authors 'Alex Wang, Amanpreet Singh, +3 authors Samuel R. Bowman' and the publication details 'Published in BlackboxNLP@EMNLP 20 April 2018 • Computer Science, Linguistics'. A TLDR section provides a brief summary of the paper's content. On the right, a 'Citations' box highlights the total of 7,460 citations, with a breakdown into 'Background Citations' (2,472), 'Methods Citations' (3,403), and 'Results Citations' (106). A 'View All' button is also present. At the bottom, a navigation bar offers links to 'Figures and Tables', 'Topics', '7,460 Citations', '77 References', and 'Related Papers'.

DOI: 10.18653/v1/W18-5446 • Corpus ID: 5034059

GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding

Alex Wang, Amanpreet Singh, +3 authors Samuel R. Bowman • Published in BlackboxNLP@EMNLP 20 April 2018 • Computer Science, Linguistics

TLDR A benchmark of nine diverse NLU tasks, an auxiliary dataset for probing models for understanding of specific linguistic phenomena, and an online platform for evaluating and comparing models, which favors models that can represent linguistic knowledge in a way that facilitates sample-efficient learning and effective knowledge-transfer across tasks.

[Expand](#)

[\[PDF\] Semantic Reader](#) [Save to Library](#) [Create Alert](#) [Cite](#)

[Figures and Tables](#) [Topics](#) [7,460 Citations](#) [77 References](#) [Related Papers](#)

7,460 Citations
Highly Influential 1,310
Citations

- Background Citations 2,472
- Methods Citations 3,403
- Results Citations 106

[View All](#)

3 Types
of Citation

The Evolution of Citation Classification Framework

12 types within 4 group

(*Weakness, Contrast, Positive, Neutral*)

(Citation Function Corpus)
CFC, 2006

6 types

(*Use, Basis, Comparison, Criticizing, Substantiating, Neutral*)

Abu-Jbara et al., 2013

3 types

Streamlined from ACL-ARC

(*Background, Method, ResultComparison*)

SciCite, 2019

CiTO, 2010

(Citation Typing Ontology)

>30 citation relations

an ontology-based
framework

ACL-ARC, 2018

6 types

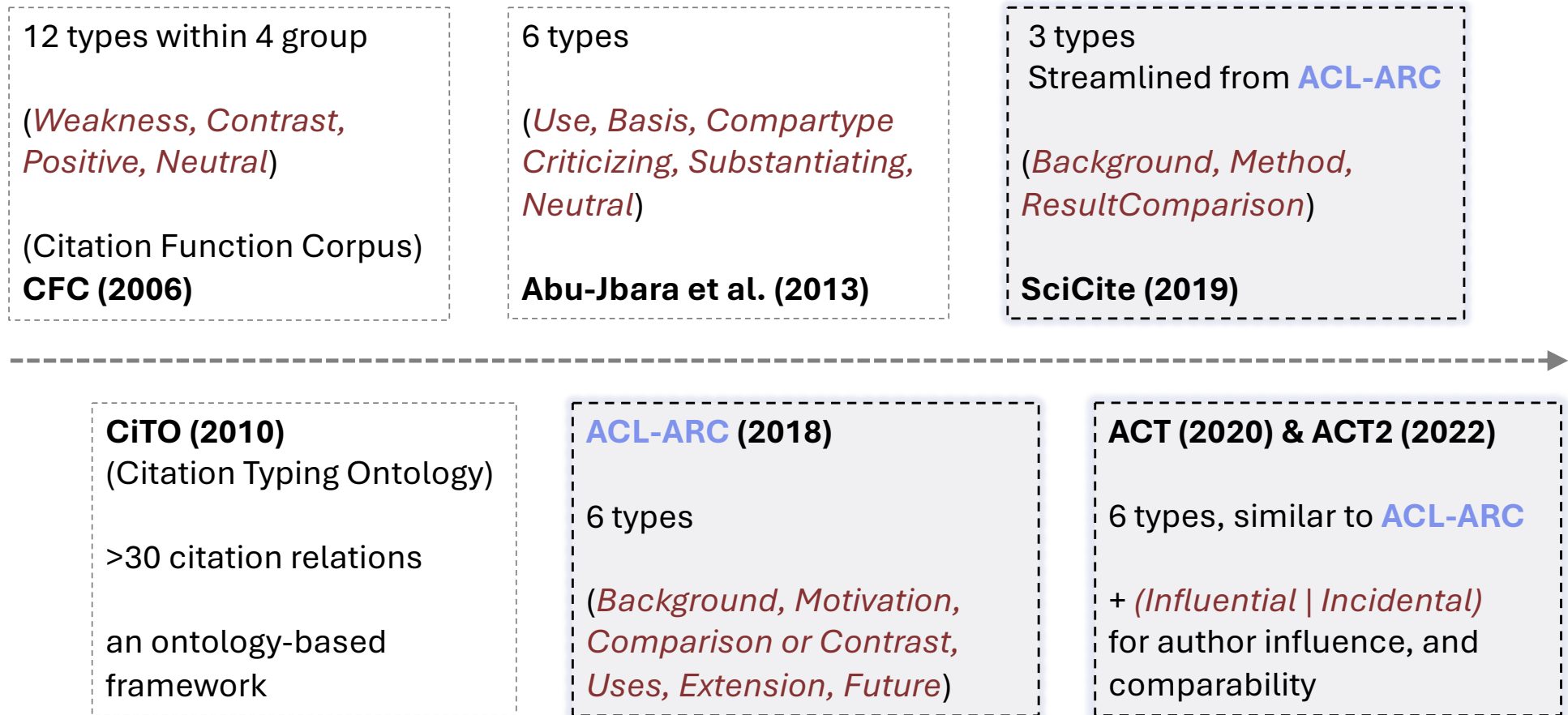
(*Background, Motivation, Comparison or Contrast, Uses, Extension, Future*)

ACT, 2020 & ACT2, 2022

6 types, similar to ACL-ARC

+ (*Influential | Incidental*)
for author influence, and
comparability

The Evolution of Citation Classification Framework



Challenges

“uses method in”,
“uses data from”,
“uses conclusions from”
—— CiTO

Forces the act of use to be
tied directly to the object
being used.

“background”,
“motivation”, “future”
(*noun-based states*)
“use”, “extend”,
“compare/contrast”
(*verb-based actions*)
—— ACL-ARC

Same grammar pattern
doesn’t fit both.
Perspective shifts between
citing vs. cited work.

“comparison/contrast”
—— ACL-ARC

“criticizing”^[2]
“substantiating”^[2]
“neutral”^[1,2]

Describe the surface act or
stance but leave the
author’s purpose implicit.

[1] Automatic classification of citation
function (Teufel et al., EMNLP 2006)

[2] Purpose and Polarity of Citation: Towards NLP-
based Bibliometrics (Abu-Jbara et al., NAACL 2013)

Challenges

Dimensional Entanglement

Conflates what is cited (method, data, finding) with why it is cited (use, extend).

This causes type explosion and leaves boundary cases unclassifiable, increasing annotator burden.

Perspective Ambiguity

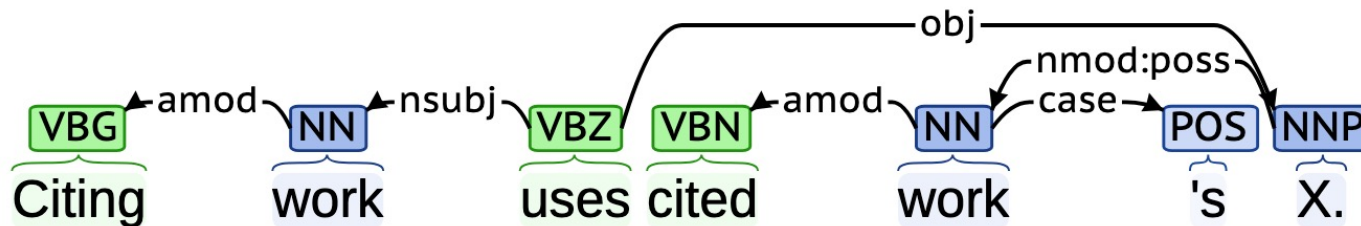
Mixes action types (verbs: use, extend, compare) with state types (nouns: background, motivation, future).

This inconsistent grammar shifts perspective between citing and cited work, confusing annotators.

Lack of Functional Grounding

describe **what** happens, but not the underlying intent.

Our Solutions



From the Perspective of **Dependency Grammar**:

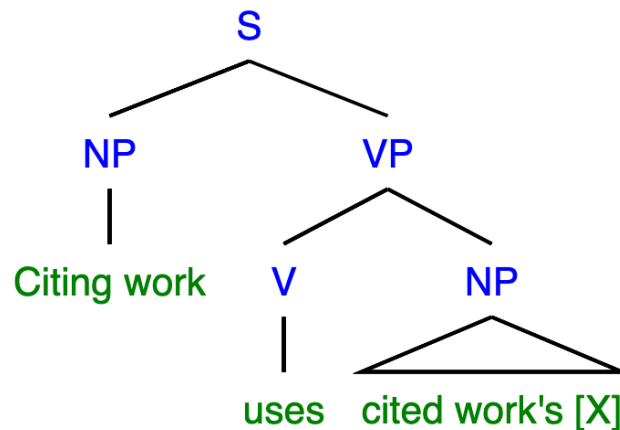
1. **Citing work** is always the **subject (nsubj)**.
2. **Citation intent** is the **root verb** of the relation.
3. **Cited work and its contribution** is always the **object (obj)**.

Disentangling
Intent and
Content

Ensuring
Perspective
consistent

(visualization with <https://corenlp.run/>)

Our Solutions



From the Perspective of **Constituency Grammar**:

1. **Cited work** and **its contribution** form a **Noun Phrase (NP)**.
2. We should also classify the types of **contributions**.

(visualization with <https://dprebyl.github.io/syntree/>)

Our Solutions – 1. Normalize Framework with Generative Grammar

S →	VP →	NP		Version 1 (ACL-ARC)
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates				Types
Citing work	frames	cited work's [X]	as background knowledge	Background
	derives	motivation	from cited work's [X]	Motivation
	uses	cited work's [X]		Uses
	extends	cited work's [X]		Extends
	compare/contrast		with cited work's [X]	Compare/Contrast
	projects	future steps	from cited work's [X]	Future

Our Solutions – 1. Normalize Framework with Generative Grammar

S →	VP →	NP	Version 1 (ACL-ARC)	
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates				Types
Citing work	frames	cited work's [X]	as background knowledge	Background
	contextualize		with cited work's [X]	Contextualize
	derives	motivation	from cited work's [X]	Motivation
	uses	cited work's [X]		Uses
	has	an extension	from cited work's [X]	Extension
	has	an comparison/contrast	with cited work's [X]	Comparison/Contrast
	projects	future steps	from cited work's [X]	Future
	justify	[Y] (future action)	with cited work's [X]	Justify

Our Solutions – 1. Normalize Framework with Generative Grammar

S →	VP →	NP	Version 1 (ACL-ARC)	
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates				Types
Citing work	frames	cited work's [X]	as background knowledge	Background
	contextualize		with cited work's [X]	Contextualize
	derives	motivation	from cited work's [X]	Motivation
	justify	[Y] (gaps/needs)	with cited work's [X]	Justify [Y]
	uses	cited work's [X]		Uses
	has	an extension	from cited work's [X]	Extension
	has	an comparison/contrast	with cited work's [X]	Comparison/Contrast
	projects	future steps	from cited work's [X]	Future

Definition of
“Motivation”
in ACL-ARC

Our Solutions – 1. Normalize Framework with Generative Grammar

S →	VP →	NP		Version 1 (ACL-ARC)
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates				Types
Citing work	frames	cited work's [X]	as background knowledge	Background
	contextualize		with cited work's [X]	Contextualize
	derives	motivation	from cited work's [X]	Motivation
	justify	[Y] (gaps/needs)	with cited work's [X]	Justify [Y]
	uses	cited work's [X]		Uses
	uses	cited work's [X]		Uses
	extends	cited work's [X]		Extends
	extends	cited work's [X]		Extends
	compare/contrast		with cited work's [X]	Compare/Contrast
	compare/contrast		with cited work's [X]	Compare/Contrast
	projects	future steps	from cited work's [X]	Future

Our Solutions – 1. Normalize Framework with Generative Grammar

S →	VP →	NP		Version 1 (ACL-ARC)
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates				Types
Citing work	frames	cited work's [X]	as background knowledge	Background
	contextualize		with cited work's [X]	Contextualize
	derives	motivation	from cited work's [X]	Motivation
	justify	[Y] (gaps/needs)	with cited work's [X]	Justify [Y]
	uses	cited work's [X]		Uses
	uses	cited work's [X]		Uses
	extends	cited work's [X]		Extends
	extends	cited work		Extends
	compare/contrast		cited work's [X]	Compare/Contrast
	compare/contrast		cited work's [X]	Compare/Contrast
	projects	future steps	from cited work's [X]	Future
	justify	[Y] (future action)	with cited work's [X]	Justify [Y]

“Future” defined as
“potential avenue”
in ACL-ARC

Our Solutions – 2. Pragmatic Distinction

S →	VP →	NP		Version 2
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates				Types
Citing work	frames	cited work's [X]	as background knowledge	Background
	contextualize		with cited work's [X]	Contextualize
	derives	motivation	from cited work's [X]	Motivation
	justify	[Y] (gaps/needs)	with cited work's [X]	Justify [Y]
	uses	cited work's [X]		Uses
	uses	cited work's [X]		Uses
	extends	cited work's [X]		Extends
	extends	cited work's [X]		Extends
	compare/contrast		with cited work's [X]	Compare/Contrast
	compare/contrast		with cited work's [X]	Compare/Contrast
	projects	future steps	from cited work's [X]	Future
	justify	[Y] (future action)	with cited work's [X]	Justify [Y]

Our Solutions – 2. Pragmatic Distinction

S →	VP →	NP	Version 2
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase
Templates			Types
Citing work	contextualize	with cited work's [X]	Contextualize
	justify	[Y] (gaps/future action) with cited work's [X]	Justify [Y]
	uses	cited work's [X]	Uses
	extends	cited work's [X]	Extends
	compare/contrast	with cited work's [X]	Compare/Contrast

Our Solutions – 2. Pragmatic Distinction

S →	VP →	NP		Version 2
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates				Types
Citing work	contextualize		with cited work's [X]	Contextualize
	justify	[Y] (gaps/future action)	with cited work's [X]	Justify
	signals	gap	from cited work's [X]	Signal Gap
	highlights	limitation	of cited work's [X]	Highlight Limitation
	justifies	design choice	with cited work's [X]	Justify Design Choice
	uses	cited work's [X]		Uses
	extends	ci		Extends
	compare/contrast		with cited work's [X]	Compare/Contrast

Gap: unresolved

Limitation: flaw or constrain of existing [X]

Design choice: design choice by referencing [X]

Our Solutions – 2. Pragmatic Distinction

S →	VP →	NP	Version 2
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase
Templates			Types
Citing work	contextualize	with cited work's [X]	Contextualize
	justify	with cited work's [X]	Justify
	signals	from cited work's [X]	
	highlights	of cited work's [X]	
	justifies	with cited work's [X]	
	uses		Uses
	extends	cited work's [X]	Extends
	modifies	cited work's [X]	Modify
	compare/contrast	with cited work's [X]	Compare/Contrast

“Modify” is broader than “Extend”: it can include adding, removing, adapting, changing, or integrating a cited work’s contribution.

Our Solutions – 2. Pragmatic Distinction

S →	VP →	NP	Version 2	
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates			Types	
Citing work	contextualize		with cited work's [X]	Contextualize
	justify	[Y] (gaps/future action)	with cited work's [X]	Justify
	signals	gap	from cited work's [X]	
	highlights	limitation	of cited work's [X]	
	justifies	design choice	with cited work's [X]	
	uses	cited work's [X]		Uses
	extends	cited work's [X]		Extends
	modifies	cited work's [X]		Modify
	compare/contrast		with cited work's [X]	Compare/Contrast
	evaluates against	cited work's [X]		Evaluates Against

Our Solutions – 2. Pragmatic Distinction

S →	VP →	NP	Version 3	
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates			Types	
Citing work	contextualize		with cited work's [X]	Contextualize
	justify	[Y] (gaps/future action)	with cited work's [X]	Justify
	signals	gap	from cited work's [X]	Signal Gap
	highlights	limitation	of cited work's [X]	Highlight Limitation
	justifies	design choice	with cited work's [X]	Justify Design Choice
	uses	cited work's [X]		Uses
	extends	cited work's [X]		Extends
	modifies	cited work's [X]		Modify
	compare/contrast		with cited work's [X]	Compare
	evaluates against	cited work's [X]		Evaluates Against

Our Solutions – 2. Pragmatic Distinction

S →	VP →	NP	Version 3	
Noun Phrase	Verb	Noun Phrase	Prep. + Noun Phrase	
Templates			Types	
Citing work	contextualize		with cited work's [X]	Contextualize
	signals	gap	from cited work's [X]	Signal Gap
	highlights	limitation	of cited work's [X]	Highlight Limitation
	justifies	design choice	with cited work's [X]	Justify Design Choice
	uses	cited work's [X]		Uses
	modifies	cited work's [X]		Modify
	evaluates against	cited work's [X]		Evaluates Against
	<div>[X] can be method, data, conclusion, or the work itself.</div>			

[X] can be method,
data, conclusion, or
the work itself.

Our Solutions – 3. Classify Cited work's [X]

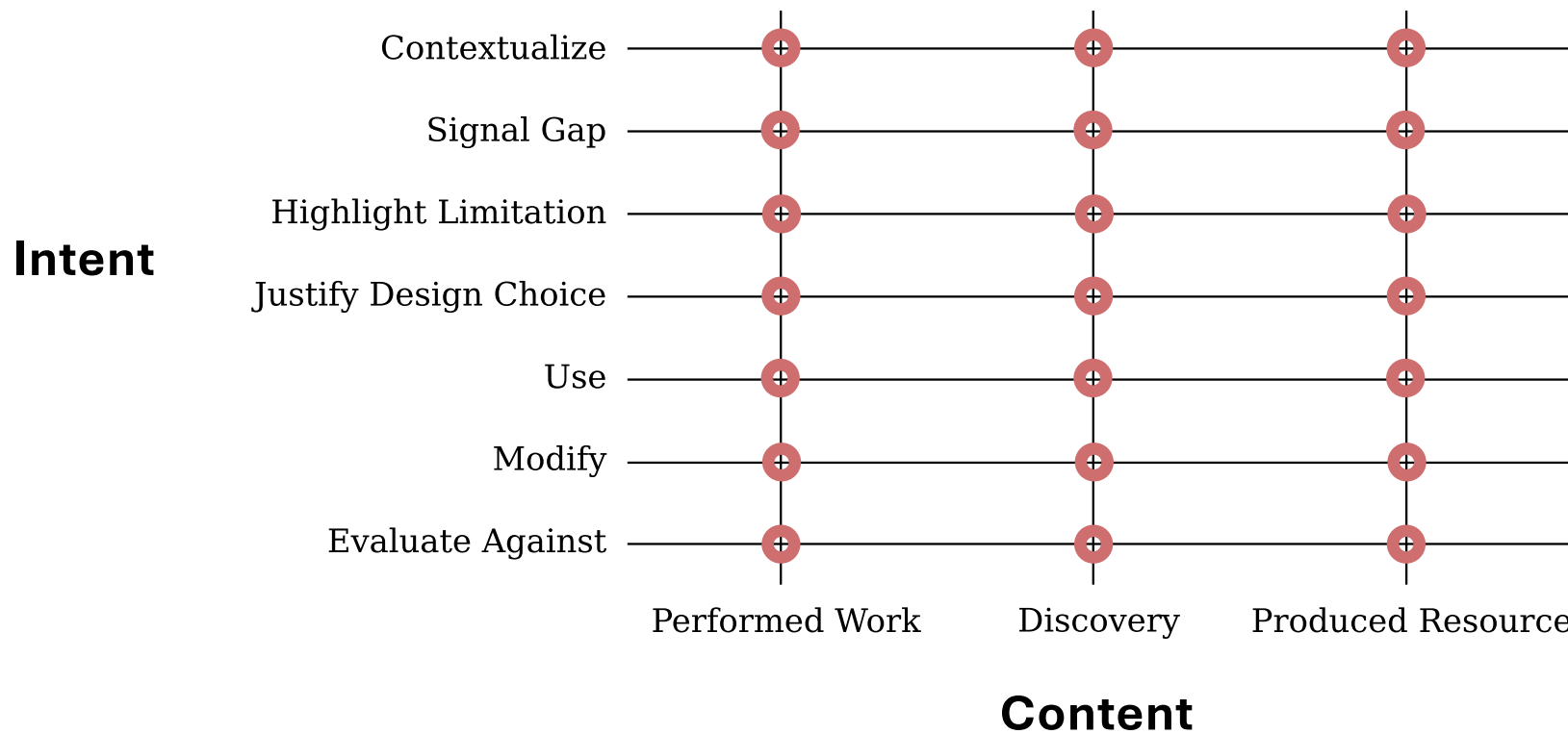
CFC (2006)	approach, methods, tools, algorithms, data, definitions	goals, results	the cited work
CiTO (2010)	method, data, solution	ideas, facts, assertion, conclusions, statements (factual evidence), information, excerpts	
Abu-Jbara et al. (2013)	approach, method, idea, tool	claims, results	the cited work
ACL-ARC (2018)	data, methods,	relevant information, goals	
SciCite (2019)	method, tool, approach, dataset	problem, concept, topic, results, findings	
ACT & ACT 2 (2020, 2022)	methodology, tools, data	information	

Our Solutions – 3. Classify Cited work's [X]

	Produced Resource	Discovery (claims)	Perform Work
CFC (2006)	approach, methods, tools, algorithms, data, definitions	goals, results	the cited work
CiTO (2010)	method, data, solution	ideas, facts, assertion, conclusions, statements (factual evidence), information, excerpts	
Abu-Jbara et al. (2013)	approach, method, idea, tool	claims, results	the cited work
ACL-ARC (2018)	data, methods,	relevant information, goals	
SciCite (2019)	method, tool, approach, dataset	problem, concept, topic, results, findings	
ACT & ACT 2 (2020, 2022)	methodology, tools, data	information	

Our Solutions – Final Framework: SOFT

Semantically **O**rthogonal **F**ramework with **T**wo dimensions



Methodology: how to evaluate the new framework SOFT?

Same Dataset: ACL-ARC (Computational Linguistics domain; Train: 1647, Test: 284)
Different Frameworks

We used preprocessed dataset from CitePrompt, which

1. filtered out bad cases that without context.
e.g. “[CITED_AUTHOR].”
2. Ensure that duplicate contexts (pointing to different citations) do not appear in both the training and test sets.

Methodology: how to evaluate the new framework SOFT?

Same Dataset: ACL-ARC (Computational Linguistics domain; Train: 1647, Test: 284)

Different Frameworks

1. Original ACL-ARC Framework

2. **Mapped** to SciCite Framework

(using the mapping method described in SciCite's paper)

Uses	→	Method
Compare / Contrast	→	ResultComparison
Background	→	Background
Extends	→	Background
Motivation	→	Background
Future	→	Background

3. **Re-annotate** with SOFT

at least three annotators for each data entry, with master degree

Methodology: how to evaluate the new framework SOFT?

Same Dataset: ACL-ARC (Computational Linguistics domain; Train: 1647, Test: 284)

Different Frameworks

An example

Citation context:

“The last years have seen considerable advances in the field of anaphora resolution, but a number of outstanding issues either remain unsolved or need more attention and, as a consequence, represent major challenges to the further development of the field ([CITED_AUTHOR]a).”

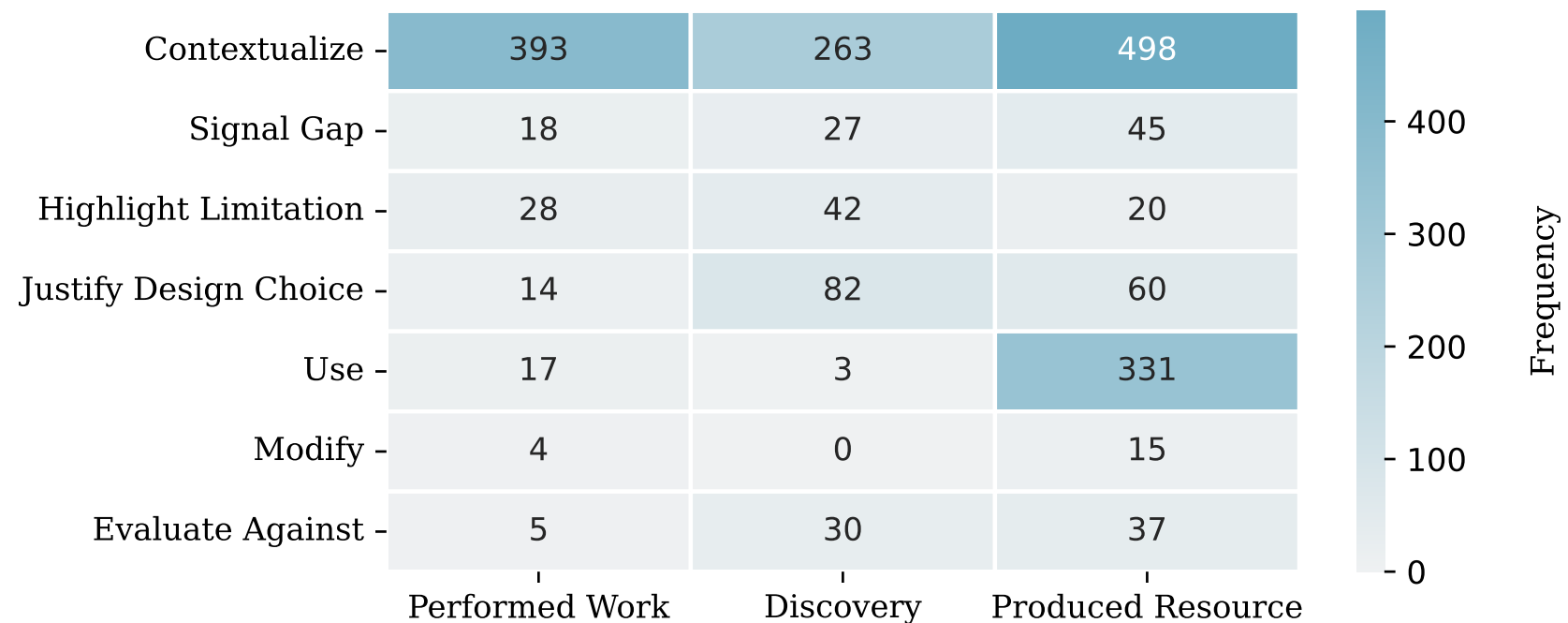
Framework	Type(s)
SciCite	Background
ACL-ARC	Motivation
SOFT (Our Framework)	Signal Gap (Intent), Discovery (Content)

Methodology: how to evaluate the new framework SOFT?

Same Dataset: ACL-ARC (Computational Linguistics domain; Train: 1647, Test: 284)

Different Frameworks

Re-annotate with SOFT - Statistics



Methodology: how to evaluate the new framework SOFT?

Same Dataset: ACL-ARC (Computational Linguistics domain; Train: 1647, Test: 284)

Different Frameworks:

Original ACL-ARC, **Mapped** SciCite, **Re-annotate** with SOFT

Same Automated Classification Method:

1. Fine-tuned Models GPU H100, 5 Hours in total
SciBert (CitePrompt, 2023), Qwen-2.5-14B (CitationIntentOpenLLM, 2025)
2. Zero-shot LLMs GPU H100, 20 mins in total
Qwen-2.5-72B, Llama-3.3-70B, Gemma-3-27B, Mistral-Small-24B

Same Evaluation Metrics: Macro F1 Score

Results

Macro F1 Score

	Classifiers	ACL-ARC 6 types	SciCite 3 types	SOFT (Content) 3 types	SOFT (Intent) 7 types
Zero-shot	Llama	0.52	0.63	0.67	0.72
	Mistral	0.49	0.59	0.62	0.71
	Gemma	0.48	0.63	0.68	0.59
	Qwen	0.45	0.59	0.67	0.75
Fine-tuned	SciBERT	0.55	0.69	0.70	0.53
	Qwen-Small	0.51	0.61	0.78	0.65
Average		0.50	0.62	0.69	0.66

1. **SOFT content** achieve **highest average** score.
 - 3 types SOFT content > 3 types SciCite
2. **SOFT Intent** retains comparable performance with 7 types.

Frameworks

Results

Macro F1 Score

	Classifiers	ACL-ARC 6 types	SciCite 3 types	SOFT (Content) 3 types	SOFT (Intent) 7 types
Zero-shot	Llama	0.52	0.63	0.67	0.72
	Mistral	0.49	0.59	0.62	0.71
	Gemma	0.48	0.63	0.68	0.59
	Qwen	0.45	0.59	0.67	0.75
Fine-tuned	SciBERT	0.55	0.69	0.70	0.53
	Qwen-Small	0.51	0.61	0.78	0.65
Average		0.50	0.62	0.69	0.66

Models

- 1. **Content types** (3) are **easier** for (fine-tuned) models to **learn** than Intent types (7).
- 2. **Zero-shot LLMs** perform **better** at handling **intent types** (7) than content types (3).
(LLMs need more domain knowledge for understanding the terminology within the content?)

Alternative Evaluations

1. **Interpretability:** evaluated by **Inter-annotation Agreement**

- Pairwise Cohen's κ : $[-1, 1]$
(No agreement: <0 , Poor: <0.2 , Fair: <0.4 , Moderate: <0.6 , Good: <0.8 , Perfect: ≤ 1)
- Not reported by previous study (not fixed annotators group for each data entry)
- Annotators: LLMs, Human
- Calculate between: LLM-LLM, Human-LLM

Alternative Evaluations

1. Interpretability: evaluated by Inter-annotation Agreement

- Pairwise Cohen's κ : $[-1, 1]$
(No agreement: <0 , Poor: <0.2 , Fair: <0.4 , Moderate: <0.6 , Good: <0.8 , Perfect: ≤ 1)

ACL-ARC (6 types)

	Llama	Mistral	Gemma	Qwen
Mistral	0.5833			
Gemma	0.4504	0.3860		
Qwen	0.5776	0.6013	0.3600	
HUMAN	0.3956	0.4048	0.4302	0.3809

SciCite (3 types)

	Llama	Mistral	Gemma	Qwen
Mistral	0.6363			
Gemma	0.6855	0.6423		
Qwen	0.6762	0.7151	0.6928	
HUMAN	0.4016	0.3850	0.4618	0.4116

SOFT Content (3 types)

	Llama	Mistral	Gemma	Qwen
Mistral	0.5396			
Gemma	0.5751	0.5657		
Qwen	0.5496	0.6169	0.5717	
HUMAN	0.5193	0.4932	0.5901	0.5596

SOFT Intent (7 types)

	Llama	Mistral	Gemma	Qwen
Mistral	0.6154			
Gemma	0.5941	0.5949		
Qwen	0.6282	0.5887	0.6175	
HUMAN	0.6620	0.6070	0.6232	0.6918

Alternative Evaluations

1. **Interpretability:** evaluated by **Inter-annotation Agreement**
2. **Generalizability:** evaluated on a cross-domain dataset

Dataset: ACT2

- 19 disciplines (cross-domain)
- Randomly select 284 data entries from the ACT2's test set
- Mapped to SciCite, re-annotate with SOFT
- Tested on zero-shot LLMs and fine-tuned models trained on the ACL-ARC dataset

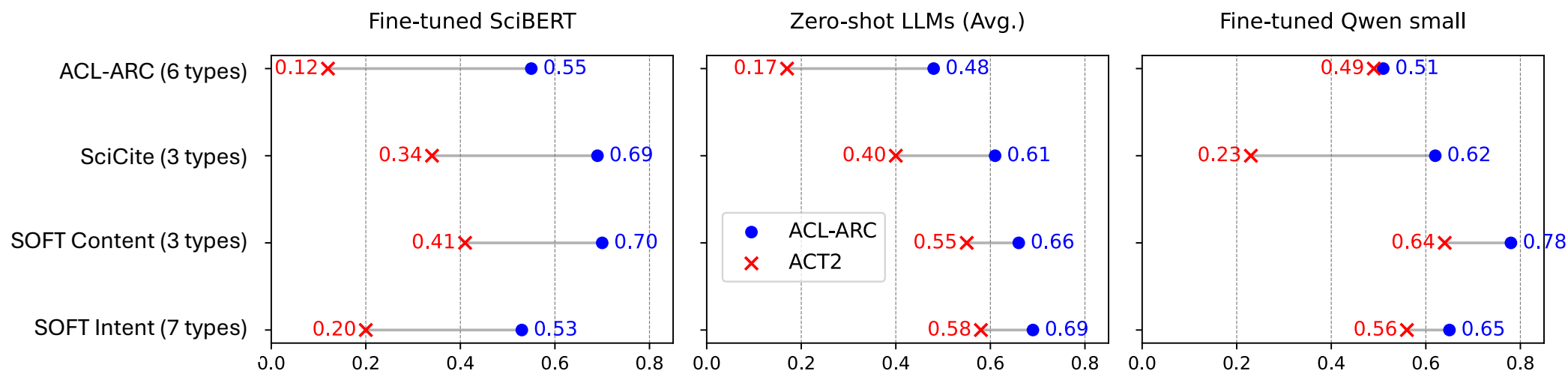
Alternative Evaluations

1. **Interpretability:** evaluated by **Inter-annotation Agreement**
2. **Generalizability:** evaluated on a cross-domain dataset

● in-domain ✕ cross-domain

highest Macro F1 Score:
(cross-domain)

- SOFT Content: 0.64 (Fine-tuned Qwen small)
- SOFT Intent: 0.58 (Zero-shot LLMs)
- ACL-ARC: 0.49 (Fine-tuned Qwen small)
- SciCite: 0.40 (Zero-shot LLMs)



Alternative Evaluations

1. **Interpretability**: evaluated by **Inter-annotation Agreement**
2. **Generalizability**: evaluated on a cross-domain dataset

✗ cross-domain - Frameworks

- Both **SOFT Content** and **SOFT Intent** drop less than the other two framework
- Also achieve **obviously higher Macro F1 Scores** than the other two framework

	Classifiers	ACL-ARC 6 types	SciCite 3 types	SOFT (Content) 3 types	SOFT (Intent) 7 types	Avg.
Zero-shot	Llama	0.13 (75%↓)	0.37 (54%↓)	0.55 (18%↓)	0.57 (21%↓)	0.41 (42%↓)
	Mistral	0.16 (67%↓)	0.41 (31%↓)	0.46 (26%↓)	0.57 (20%↓)	0.40 (36%↓)
	Gemma	0.17 (65%↓)	0.36 (43%↓)	0.58 (15%↓)	0.55 (7%↓)	0.42 (33%↓)
	Qwen	0.23 (49%↓)	0.45 (24%↓)	0.60 (10%↓)	0.64 (15%↓)	0.48 (25%↓)
Fine-tuned	SciBERT	0.12 (78%↓)	0.34 (51%↓)	0.41 (41%↓)	0.20 (62%↓)	0.27 (58%↓)
	Qwen-Small	0.49 (4%↓)	0.23 (63%↓)	0.64 (18%↓)	0.56 (14%↓)	0.48 (25%↓)
	Avg.	0.22 (56%↓)	0.36 (44%↓)	0.54 (21%↓)	0.52 (23%↓)	

Alternative Evaluations

1. **Interpretability**: evaluated by **Inter-annotation Agreement**
2. **Generalizability**: evaluated on a cross-domain dataset

✗ cross-domain - Models

- Both zero-shot & fine-tuned **Qwen** achieved best performance
- **SciBERT** performed worst in the cross-domain task

	Classifiers	ACL-ARC 6 types	SciCite 3 types	SOFT (Content) 3 types	SOFT (Intent) 7 types	Avg.
Zero-shot	Llama	0.13 (75%↓)	0.37 (54%↓)	0.55 (18%↓)	0.57 (21%↓)	0.41 (42%↓)
	Mistral	0.16 (67%↓)	0.41 (31%↓)	0.46 (26%↓)	0.57 (20%↓)	0.40 (36%↓)
	Gemma	0.17 (65%↓)	0.36 (43%↓)	0.58 (15%↓)	0.55 (7%↓)	0.42 (33%↓)
	Qwen	0.23 (49%↓)	0.45 (24%↓)	0.60 (10%↓)	0.64 (15%↓)	0.48 (25%↓)
Fine-tuned	SciBERT	0.12 (78%↓)	0.34 (51%↓)	0.41 (41%↓)	0.20 (62%↓)	0.27 (58%↓)
	Qwen-Small	0.49 (4%↓)	0.23 (63%↓)	0.64 (18%↓)	0.56 (14%↓)	0.48 (25%↓)
	Avg.	0.22 (56%↓)	0.36 (44%↓)	0.54 (21%↓)	0.52 (23%↓)	