

Grounding “grounding”: How has grounding evolved in meaning?

Eric Huang

University of Waterloo
e48huang@uwaterloo.ca

Abstract

Terminology used within linguistics and AI conferences tend to be overused, leading to ambiguous meaning and difficulty navigating new papers. This paper will elucidate the various senses of the word “grounding” through qualitative analysis and deeper quantitative analysis of its various uses. In particular, we aim to show how the senses of the word has evolved throughout the years of various conferences. All code can be found at https://git.uwaterloo.ca/e48huang/cs-784/-/tree/final_project/final_project?ref_type=heads where a University of Waterloo account is required.

Conference	Paper Count
AAAI	772
ACL	632
CVPR	862
ECCV	511
EMNLP	575
ICCV	341
ICLR	360
ICML	360
IJCAI	654
NAACL	226
NeurIPS	654

Table 1: Counts of unique papers with “grounding” by conference found in the corpora.

1 Introduction

Many conferences centering around Artificial Intelligence have existed for many decades, evolving over time on the types of problems that they tackle. While these problems change over time, so do the terminology, which have a tendency to evolve semantically, leading to overloaded terms. One such term is “grounding”, the idea that one wishes to ensure that there is understanding or a common ground (Nakano et al., 2003). While this term seems simple, it is used in many various contexts, all of which requires different datasets, methods and metrics to evaluate, while being applied in different settings.

To better understand the term “grounding” and its usage, we perform both quantitative analysis and qualitative analysis. This paper explores the “Seed42Lab/AI-paper-crawl” HuggingFace dataset (Forty-Two AI Lab) which collects full-text papers from 11 different conferences spanning from the first year of the conference to 2024. To first select different senses of the word “grounding”, we perform preliminary quantitative analysis to filter for papers to further investigate. From these selected papers, we identify 8 related but distinct meanings

of the word “grounding”. We perform some literature review to understand how these different senses are understood, from its various datasets, methods, metrics and applications. Finally, for each of these word senses, we investigate how they have evolved over time.

2 Paper Selection

A simple search over the number of papers which have the term “grounding” quickly shows that it is infeasible to cover all possible instances. For example, the Association for Computational Linguistics (ACL) alone has 632 unique papers that have an instance of “grounding” (see Table 1). While not all these instances are due to the paper itself being related to grounding, as they can simply include a citation within its bibliography, they are indicative that some filtering of papers is necessary.

To filter through these papers, we propose a method which selects the most relevant papers within a conference to the word “grounding”. We take a naive approach where we select the top 10% of papers with the word “grounding”. We determine which papers are more important to “ground-

ing” based on the word frequency, if “grounding” appears more often compared to other words within a paper, then it should be more relevant. This process (while ensuring uniqueness across conferences) resulted in 46 curated papers¹, spanning from the years 2000 to 2024, shown in the following figure. While this selection of papers may not cover the breadth of senses that grounding might entail, as it misses on papers from the 1980’s to 2000’s, it does cover the most commonly used senses of the word.

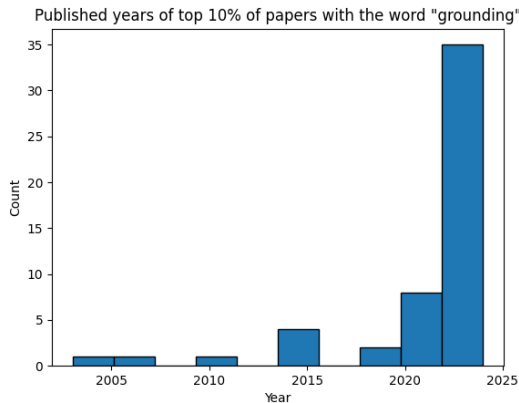


Figure 1: Count of selected conference papers per year

3 Grounding “grounding”

In this section we cover the 8 different senses of the word “grounding” found in the 46 papers covered in the previous section. We will explore the meanings of each sense and the challenges that each paper is tackling. The delineation of these different senses tend to be ambiguous; understanding whether certain senses belong to certain categories can be argued. Rather, categories were chosen according to modality and how differently “grounding” would be utilized and established.

3.1 Visual Grounding

Visual grounding, also known as image, phrase or referring expression grounding (Xiao et al., 2024; Li et al., 2024; Ma et al., 2020; Islam et al., 2023; Jiang et al., 2019; Lu et al., 2022; Dou and Peng, 2021; Surikuchi et al., 2023) refers to the challenge of trying to localize specific regions within an image based on some textual description². Traditionally,

¹https://git.uwaterloo.ca/e48huang/cs-784/-/blob/e09a1c22c0de7e331ca16109a5f32b226dc6d9c5/final_project/grounding_top_p.txt

²other terms include natural language object retrieval or phrase localization (Ma et al., 2024)

this involved finding the phrase’s referring region by predicting a bounding box around said region. As time has gone on however, there has been more and more types of challenges that one could tackle within visual grounding (Xiao et al., 2024). From the 46 papers we filtered for, we have discovered the following subcategories involved in visual grounding.

3.1.1 Classical Visual Grounding

This entails the traditional problem of trying to predict a bounding box around a region (Li et al., 2024; Huang et al., 2021; Peng et al., 2023). Recent papers have found different methods in an attempt to improve performance. Zeng et al. (2024) improves compositional understanding through providing a harder dataset that relies on introducing different compositions of objects. Zhang et al. (2020) improves performance in weakly supervised (no bounding box annotations) settings through contrastive learning. Ma et al. (2024) provides a new dataset for higher resolution images at different granularity and bounding box sizes. Lee and Sung (2024) has shown improvements in image generation as well.

This visual grounding task can be understood as the inverse problem to image captioning, where one is given an image and need to provide the text portion. In fact, this inverse paradigm has led to better models (Wang et al., 2023a) involving cyclic updates.

3.1.2 Answer Grounding

Rather than fit a bounding box to various objects, Visual Question Answering (VQA) grounding attempts to find specific parts of an image that corresponds with inputted questions rather than descriptive prompts (Chen et al., 2022, 2023).

3.2 Action Grounding

Relying on other types of grounding such as image grounding, action grounding is a term that refers to building a model that is able to take some grounding and relate it to a set of actions. Recent works utilize LLMs in the fields of chat agents, web agents and robotics (Zhang et al., 2023; Cheng et al., 2024; Zheng et al., 2024; Tellex et al., 2011; Wang et al., 2023b) to motivate better actions that are aligned with people’s understanding of the world.

Contrary to using image grounding, Kameko et al. (2015) matches certain states of games to commentary in an attempt to understand how various

actions are grounded in language or its symbols. They refer to this type of grounding as “symbol grounding” but essentially attempts to relate some action to some other observation.

3.3 Audio Grounding

Audio grounding is the task of taking static images and sounds and attempting to identify which parts of the image are correlated with certain parts of audio. For example, [Tian et al. \(2021\)](#) attempts to separate images of bands into which instruments produce what kinds of audio.

3.4 Video Grounding

Another related grounding task to visual grounding is the idea of video grounding or spatio-temporal grounding. This task is to identify various portions of a video or the entities within them to provide an understanding for a certain prompt ([Jiang et al., 2024](#)). These different grounding tasks can be split into its own categories defined in the next sections.

3.4.1 Object Tracking

Object tracking relies on the idea that given some natural language prompt, to both identify the specific object within the video but also to continuously track it throughout the video or still frames ([Zhou et al., 2023](#)).

3.4.2 Natural Language Spatial Video Grounding

This video grounding task is an extension of classic visual grounding, where the model attempts to set a bounding box for each frame of a video ([Li et al., 2022](#); [Ma et al., 2020](#)).

3.4.3 Temporal Video Grounding

This video grounding task is to identify the timestamps in which a prompt holds true for a video ([Li et al., 2024](#); [Afouras et al., 2023](#); [Bao et al., 2021](#); [Chen et al., 2018](#)).

3.4.4 Spatio-temporal Video Grounding

This video grounding task combines the last two tasks and attempts to identify both the bounding boxes and the timestamps in which a prompt holds true for a video ([Wasim et al., 2024](#); [Chen et al., 2024](#); [Jin et al., 2022](#)). It can be used within various settings including video entailment which determines whether a prompt holds true for some video ([Chen and Kong, 2021](#)). Similar to image grounding, video grounding can also be used within video generation tools ([Jeong and Ye, 2024](#)).

3.5 3D Grounding

Similar to image grounding, 3D grounding adds a dimension and attempts to put bounding boxes around 3D models which are often represented as point clouds. These 3D grounding tasks share similar strategies to image grounding, using captioning tasks to improve performance ([Cai et al., 2022](#); [Yang et al., 2023](#); [Miyanishi et al., 2023](#); [Wang et al., 2023c](#)). Some papers have even used 2D object representations to improve 3D grounding ([Yang et al., 2021](#)), while others have improved 3D visual grounding with reasoning ([Zhu et al., 2024](#)).

3.6 Dialogue Grounding

This term of “dialogue grounding” is loosely defined, usually seen in literature simply as “grounding”. Within these papers, “grounding” refers to the idea of trying to build a common ground of understanding between two or more actors within a conversation. It includes attempting to analyze nonverbal behaviours ([Nakano et al., 2003](#); [Roque, 2007](#); [Liu et al., 2012](#); [Shaikh et al., 2024](#)).

3.7 Markov Logic Networks Grounding

“Grounding” in Markov Logic Networks (MLNs) differs significantly from the other senses of the word ([Venugopal and Gogate, 2014](#)). MLNs refer to a statistical model for probabilistic logic reasoning, where by developing a set of first-order logic rules known as “grounds” one is able to form a weighted satisfiability problem with an optimized solution. In particular, grounding within Markov Logic Networks refers to the process of forming the weighted satisfiability graph ([Fang et al., 2023](#)).

3.8 Physical Dynamics Grounding

Attempting to model physical dynamics purely from states and its transitions tend to be difficult, requiring a ton of resources to supervise consecutive particle properties. Instead of requiring this supervision, a new field has emerged to attempt to understand these physical dynamics from visual observations ([Cao et al., 2024](#)). One such application is in fluid dynamics grounding; which attempts to build an understanding of fluid particle systems from sequential visual observations ([Guan et al., 2022](#)).

4 Analyzing “grounding”’s Usage

In this section, we will build a quantitative understanding of “grounding” and its senses over

time. We will explore how the word has been used throughout the years, and dive deeper into a few senses of the word. We will accomplish this through observing the co-occurrence trends over time with other key words for each specific sense.

4.1 “grounding” Over Time

In this section, we explore how the term “grounding” has evolved over time through analyzing how many papers have included the term “grounding”. We aggregate over all the data splits while showcasing a more fine-grained example for a specific conference to avoid any patterns lost through aggregation.

In particular, we observe that the number of instances of “grounding” has increased both in terms of pure count and frequency over time (see Fig 2 and Fig 3). We normalize because the number of papers being published in general increases as well, naturally inflating the number of “grounding” papers. However we observe that both the pure count and frequency increase over time, concluding that “grounding” has been a terminology that is becoming more and more utilized. This is likely due to it becoming more relevant with the uprise of multimodal models (Xiao et al., 2024) and a need to interpret and improve these models.

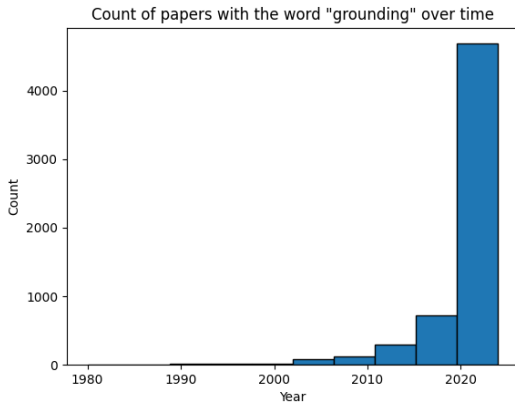


Figure 2: Count of “grounding” for all the conferences over time

To see how each individual conference’s count and frequency changes over time, see Appendix A. These graphs do confirm that our selection of papers in the previous section were well-justified, as the most important papers relevant to “grounding” are likely to be the more recent papers. Therefore, not having covered senses of “grounding” from papers spanning the 1980’s-2000’s is not as significant as it may seem.

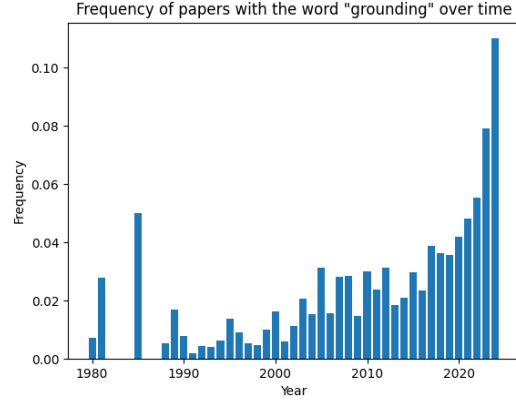


Figure 3: Frequency of papers containing the word “grounding” for all conferences over time

4.2 “grounding” Senses Over Time

We observed that the different senses of the word “grounding” are not uniformly distributed across time, rather that it has evolved. Take “dialogue grounding”, where our filtered papers were some of the only ones from the 2000s (Nakano et al., 2003; Roque, 2007) with more recent papers covering other senses of the word. To better understand “grounding”’s evolution, this section here covers different word co-occurrences over time.

In particular, we explore the frequency of papers that include the word “grounding” which also contains other words which can indicate different senses. Table 2 in Appendix B shows which words we count as co-occurring for each sense. In the following sections we explore how these co-occurrences change over time.

4.2.1 Visual Grounding

For visual grounding, we can tell that only more recently has there been an increase in the number of papers dealing with the visual grounding paradigm. According to Xiao et al. (2024), this is likely due to improvements in multimodal models in 2021, correlating with our findings in Fig 4 and Fig 5.

See Appendix C.1 for the splits per conference.

4.2.2 Action Grounding

For action grounding, there has been a steady interest over time shown by the frequency of papers which sit around 50%. At first glance, this seems high and likely to be conflated due to search words such as “web” or “agent”. However, as action grounding refers to the applicability of other types of grounding this is likely representative of the word sense itself.

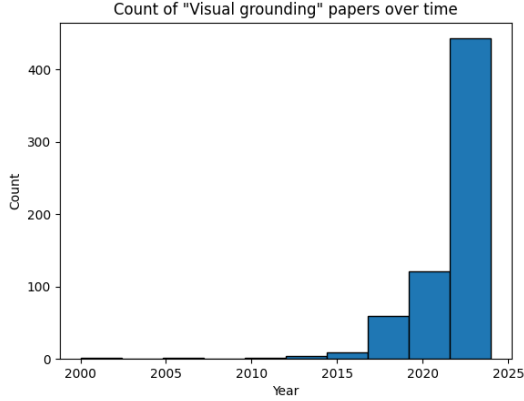


Figure 4: Count of “visual grounding” for all the conferences over time

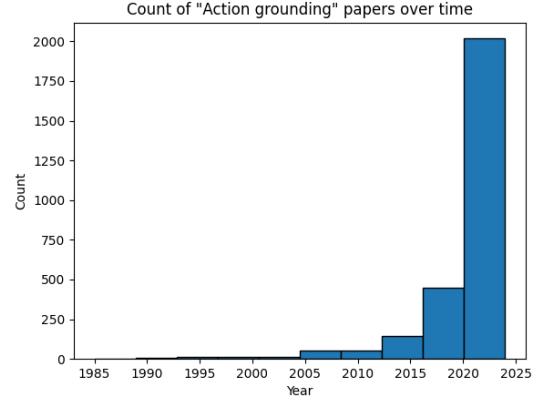


Figure 6: Count of “Action grounding” for all the conferences over time

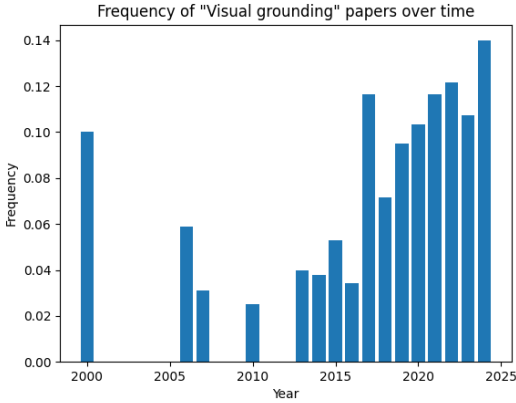


Figure 5: Frequency of “visual grounding” for all conferences over time

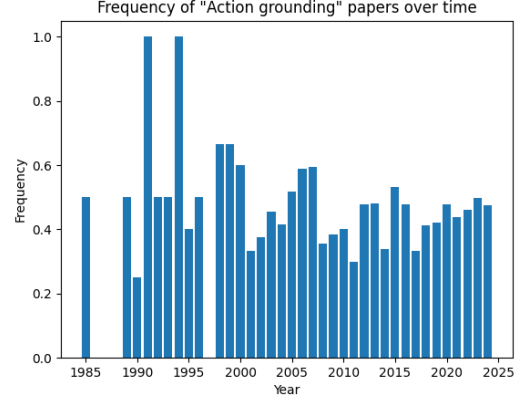


Figure 7: Frequency of “Action grounding” for all conferences over time

A notable observation is that the frequency of “grounding” papers in earlier years such as within the 1990s tend to have very high frequency that matches with “action grounding”. This is likely due to the small sample size within those time periods, having only AAAI, ACL, IJCAI and NeurIPS as conferences, each with a small magnitude of publications. This reduces the variability and thus we would expect to see higher frequencies of certain senses of words during these time periods. This trend follows for the other senses.

See Appendix C.2 for the splits per conference.

4.2.3 Audio Grounding

There seems to be less papers revolved around audio grounding, as after the 2000s, it is at most referenced in about 20% of the papers. Even as the multimodal model mark in 2021 hits, there has been a slight increase but still smaller share of the “grounding” papers count.

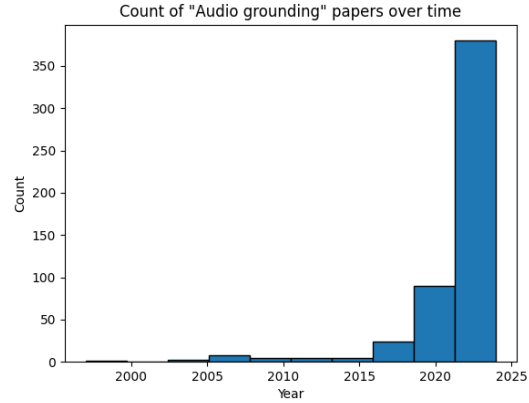


Figure 8: Count of “Audio grounding” for all the conferences over time

See Appendix C.3 for the splits per conference.

4.2.4 Video Grounding

After the 2000s, the video grounding sense follows a very similar trend to audio grounding, but with

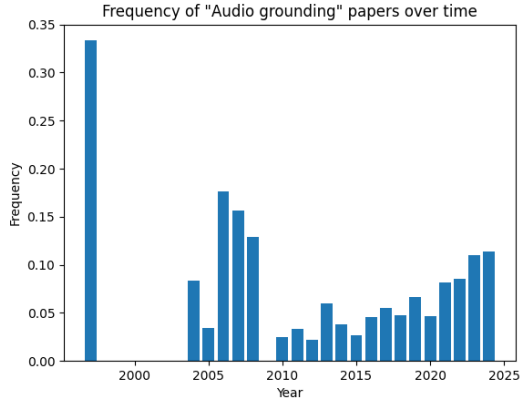


Figure 9: Frequency of “Audio grounding” for all conferences over time

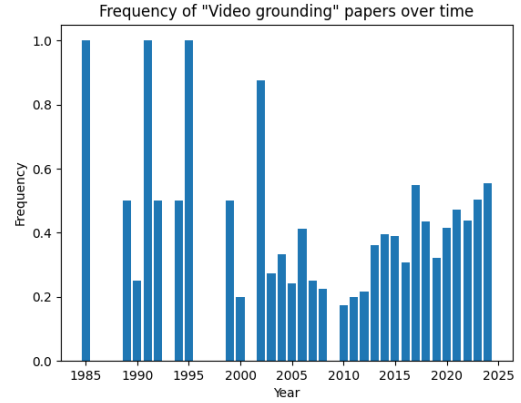


Figure 11: Frequency of “Video grounding” for all conferences over time

a significantly higher share at around 50-60% of papers. This is likely due to the fact that most multimodal models that work on video also work on other senses of grounding such as audio and image grounding (Li et al., 2024).

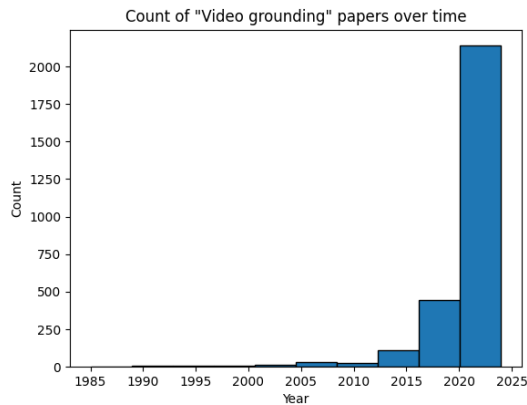


Figure 10: Count of “Video grounding” for all the conferences over time

See Appendix C.4 for the splits per conference.

4.2.5 3D Grounding

3D grounding observes a huge spike in papers around the 2020s, likely due to significant advancements in marquee papers such as ScanRefer (Chen et al., 2020; Liu et al., 2024). Such papers introduce novel problems which encourages future development and a larger share of the paper frequencies.

See Appendix C.5 for the splits per conference.

4.2.6 Dialogue Grounding

Dialogue grounding’s trend follows our empirical observations, where they had a much larger share

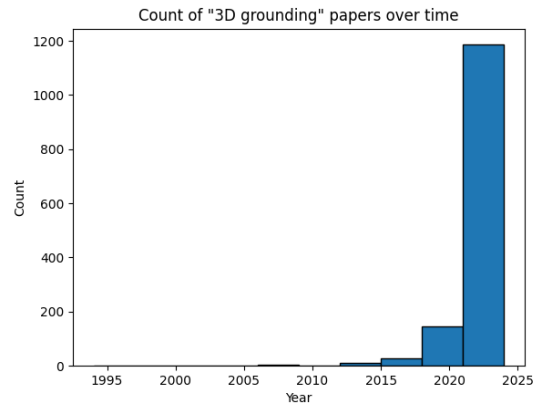


Figure 12: Count of “3D grounding” for all the conferences over time

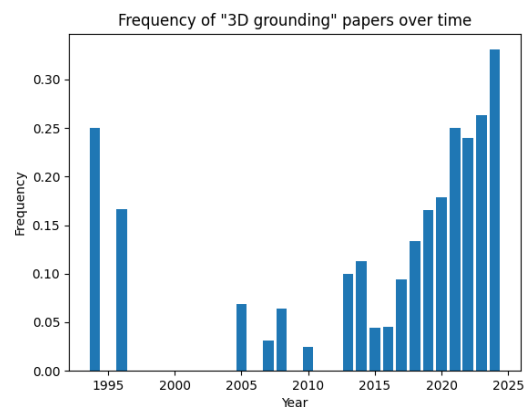


Figure 13: Frequency of “3D grounding” for all conferences over time

of papers earlier on in the 2000s, but has since decreased significantly.

See Appendix C.6 for the splits per conference.

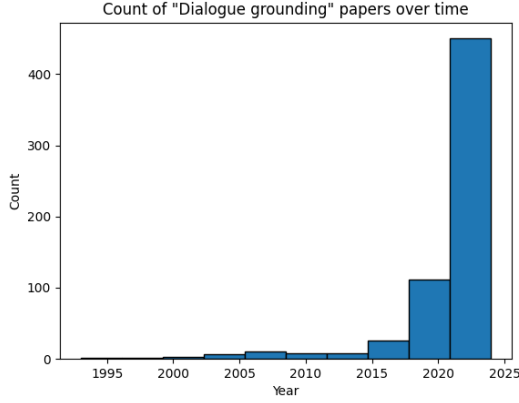


Figure 14: Count of “Dialogue grounding” for all the conferences over time

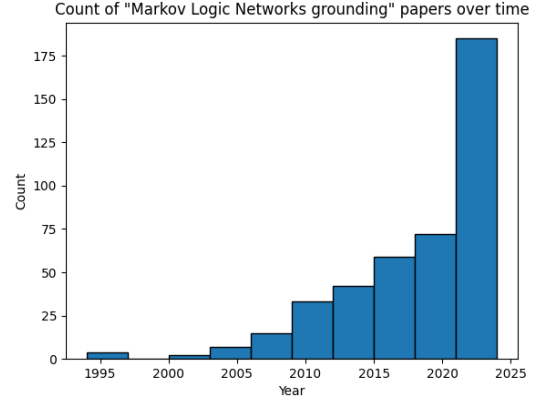


Figure 16: Count of “Markov Logic Networks grounding” for all the conferences over time

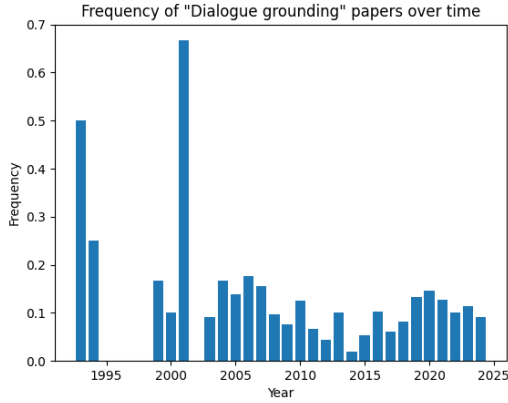


Figure 15: Frequency of “Dialogue grounding” for all conferences over time

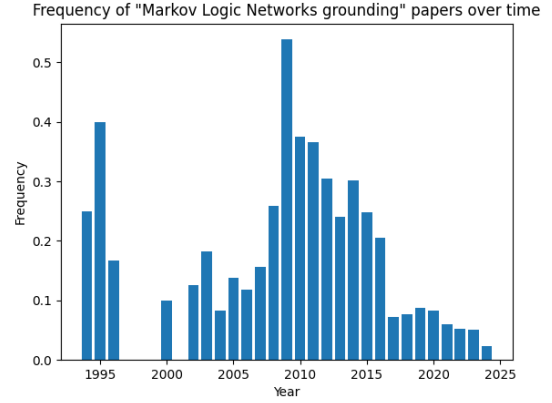


Figure 17: Frequency of “Markov Logic Networks grounding” for all conferences over time

4.2.7 Markov Logic Networks Grounding

Markov Logic Networks were seemingly popular within the 2010s, having a high share of the market at that time. However, as 2020s approached, there seems to be a shift away from Markov probabilistic models and more towards LLMs and multimodal models.

See Appendix C.7 for the splits per conference.

4.2.8 Physical Dynamics Grounding

Physical dynamics models tend to be quite niche, leading to a very small share of the amount of papers which include that sense of the word.

See Appendix C.8 for the splits per conference.

5 Conclusion

This work has shown that “grounding” is an overloaded term with many different senses and modalities. These senses have evolved over time, inflating and contracting according to research trends.

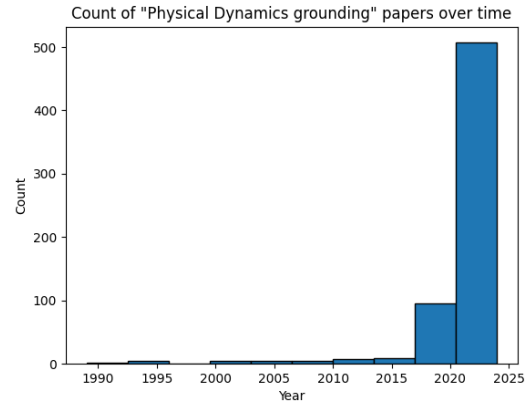


Figure 18: Count of “Physical Dynamics grounding” for all the conferences over time

In future work, we hope to explore more meta-analysis through quantitative analysis of the different datasets, methods and metrics of each sense. Especially as deep learning and multimodal models become more and more popular, understanding

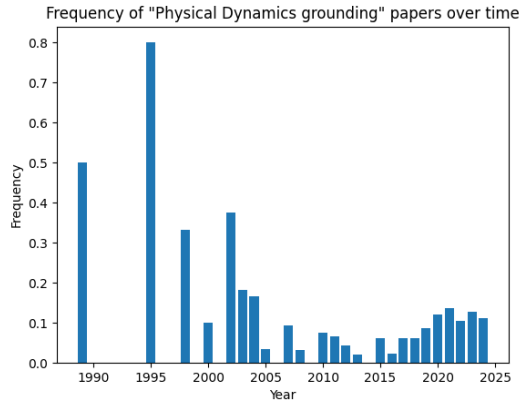


Figure 19: Frequency of “Physical Dynamics grounding” for all conferences over time

which of these datasets are being the most benchmarked could provide some more insight than a typical survey paper. We also hope to perform more fine-grained meta-analysis, understanding how much of each sense is composed of smaller sub-categories, such as understanding how much of “action grounding” is for robotics applications. Furthermore, certain word senses are also ambiguous for its high frequency and nature of encompassing other types of grounding such as “video grounding”. Future analysis is required to provide a deeper understanding of these dynamics.

Limitations

This section discusses the possible limitations from our selection process and meta-analysis. Regarding our selection process, there is potential bias in choosing the most important 10% of papers due to the nature of some work being longer with more citations and thus requiring more “grounding” occurrences to rank as important. This might lead to some over-representation of senses and a lack of other senses within the filtered papers. Similarly, the most 10% of papers might be confounded by time, as time goes on, there is an increase in the number of papers and conferences, which might artificially inflate the number of papers with more modern senses.

For our meta-analysis, our selection method of the words chosen might introduce some unwanted bias. In particular, words such as “sound” has an ambiguous semantic meaning, either referring to a sound argument or the production of noise. Future works should reduce this limitation by having better filtering in place.

References

- Triantafyllos Afouras, Effrosyni Mavroudi, Tushar Nagarajan, Huiyu Wang, and Lorenzo Torresani. 2023. [Ht-step: Aligning instructional articles with how-to videos](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 50310–50326. Curran Associates, Inc.
- Peijun Bao, Qian Zheng, and Yadong Mu. 2021. [Dense events grounding in video](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(2):920–928.
- Daigang Cai, Lichen Zhao, Jing Zhang, Lu Sheng, and Dong Xu. 2022. 3djcg: A unified framework for joint dense captioning and visual grounding on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16464–16473.
- Junyi Cao, Shanyan Guan, Yanhao Ge, Wei Li, Xiaokang Yang, and Chao Ma. 2024. [Neuma: Neural material adaptor for visual grounding of intrinsic dynamics](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 65643–65669. Curran Associates, Inc.
- Brian Chen, Nina Shvetsova, Andrew Rouditchenko, Daniel Kondermann, Samuel Thomas, Shih-Fu Chang, Rogerio Feris, James Glass, and Hilde Kuehne. 2024. What when and where? self-supervised spatio-temporal grounding in untrimmed multi-action videos from narrated instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18419–18429.
- Chongyan Chen, Samreen Anjum, and Danna Gurari. 2022. [Grounding answers for visual questions asked by visually impaired people](#). *Preprint*, arXiv:2202.01993.
- Chongyan Chen, Samreen Anjum, and Danna Gurari. 2023. Vqa therapy: Exploring answer differences by visually grounding answers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15315–15325.
- Dave Zhenyu Chen, Angel X. Chang, and Matthias Nießner. 2020. [Scanrefer: 3d object localization in rgb-d scans using natural language](#). *Preprint*, arXiv:1912.08830.
- Jingyuan Chen, Xinpeng Chen, Lin Ma, Zequn Jie, and Tat-Seng Chua. 2018. [Temporally grounding natural sentence in video](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 162–171, Brussels, Belgium. Association for Computational Linguistics.
- Junwen Chen and Yu Kong. 2021. Explainable video entailment with grounded visual evidence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2021–2030.

- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Li YanTao, Jianbing Zhang, and Zhiyong Wu. 2024. [SeeClick: Harnessing GUI grounding for advanced visual GUI agents](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9313–9332, Bangkok, Thailand. Association for Computational Linguistics.
- Zi-Yi Dou and Nanyun Peng. 2021. [Improving pre-trained vision-and-language embeddings for phrase grounding](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6362–6371, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Huang Fang, Yang Liu, Yunfeng Cai, and Mingming Sun. 2023. [Mln4kb: an efficient markov logic network engine for large-scale knowledge bases and structured logic rules](#). In *Proceedings of the ACM Web Conference 2023, WWW '23*, page 2423–2432, New York, NY, USA. Association for Computing Machinery.
- Forty-Two AI Lab. Seed42lab/ai-paper-crawl. <https://huggingface.co/datasets/Seed42Lab/AI-paper-crawl>.
- Shanyan Guan, Huayu Deng, Yunbo Wang, and Xiaokang Yang. 2022. [Neurofluid: Fluid dynamics grounding with particle-driven neural radiance fields](#). *Preprint*, arXiv:2203.01762.
- Jianqiang Huang, Yu Qin, Jiaxin Qi, Qianru Sun, and Hanwang Zhang. 2021. [Deconfounded visual grounding](#). *Preprint*, arXiv:2112.15324.
- Md Mofijul Islam, Alexi Gladstone, and Tariq Iqbal. 2023. [Patron: Perspective-aware multitask model for referring expression grounding using embodied multimodal cues](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(1):971–979.
- Hyeonho Jeong and Jong Chul Ye. 2024. [Ground-a-video: Zero-shot grounded video editing using text-to-image diffusion models](#). *Preprint*, arXiv:2310.01107.
- Ming Jiang, Qiuyuan Huang, Lei Zhang, Xin Wang, Pengchuan Zhang, Zhe Gan, Jana Diesner, and Jianfeng Gao. 2019. [TIGER: Text-to-image grounding for image caption evaluation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2141–2152, Hong Kong, China. Association for Computational Linguistics.
- Wenhui Jiang, Yibo Cheng, Linxin Liu, Yuming Fang, Yuxin Peng, and Yang Liu. 2024. [Comprehensive visual grounding for video description](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(3):2552–2560.
- Yang Jin, Yongzhi Li, Zehuan Yuan, and Yadong Mu. 2022. [Embracing consistency: A one-stage approach for spatio-temporal video grounding](#). *Preprint*, arXiv:2209.13306.
- Hiroataka Kameko, Shinsuke Mori, and Yoshimasa Tsu-ruoka. 2015. [Can symbol grounding improve low-level NLP? word segmentation as a case study](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2298–2303, Lisbon, Portugal. Association for Computational Linguistics.
- Phillip Y. Lee and Minhyuk Sung. 2024. [Reground: Improving textual and spatial grounding at no cost](#). *Preprint*, arXiv:2403.13589.
- Mengze Li, Tianbao Wang, Haoyu Zhang, Shengyu Zhang, Zhou Zhao, Jiaxu Miao, Wenqiao Zhang, Wenming Tan, Jin Wang, Peng Wang, Shiliang Pu, and Fei Wu. 2022. [End-to-end modeling via information tree for one-shot natural language spatial video grounding](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8707–8717, Dublin, Ireland. Association for Computational Linguistics.
- Zhaowei Li, Qi Xu, Dong Zhang, Hang Song, YiQing Cai, Qi Qi, Ran Zhou, Junting Pan, Zefeng Li, Vu Tu, Zhida Huang, and Tao Wang. 2024. [GroundingGPT: Language enhanced multi-modal grounding model](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6657–6678, Bangkok, Thailand. Association for Computational Linguistics.
- Changsong Liu, Rui Fang, and Joyce Chai. 2012. [Towards mediating shared perceptual basis in situated dialogue](#). In *Proceedings of the 13th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 140–149, Seoul, South Korea. Association for Computational Linguistics.
- Daizong Liu, Yang Liu, Wencan Huang, and Wei Hu. 2024. [A survey on text-guided 3d visual grounding: Elements, recent advances, and future directions](#). *Preprint*, arXiv:2406.05785.
- Panzhong Lu, Xin Zhang, Meishan Zhang, and Min Zhang. 2022. [Extending phrase grounding with pronouns in visual dialogues](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7614–7625, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Chih-Yao Ma, Yannis Kalantidis, Ghassan AlRegib, Peter Vajda, Marcus Rohrbach, and Zsolt Kira. 2020. [Learning to generate grounded visual captions without localization supervision](#). *Preprint*, arXiv:1906.00283.
- Tao Ma, Bing Bai, Haozhe Lin, Heyuan Wang, Yu Wang, Lin Luo, and Lu Fang. 2024. When visual grounding meets gigapixel-level large-scale scenes: Benchmark

- and approach. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22119–22128.
- Taiki Miyanishi, Fumiya Kitamori, Shuhei Kurita, Jungdae Lee, Motoaki Kawanabe, and Nakamasa Inoue. 2023. [Cityrefer: Geography-aware 3d visual grounding dataset on city-scale point cloud data](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 77758–77770. Curran Associates, Inc.
- Yukiko Nakano, Gabe Reinstein, Tom Stocky, and Justine Cassell. 2003. [Towards a model of face-to-face grounding](#). In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 553–561, Sapporo, Japan. Association for Computational Linguistics.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. [Kosmos-2: Grounding multimodal large language models to the world](#). *Preprint*, arXiv:2306.14824.
- Antonio Roque. 2007. [Reacting to agreement and error in spoken dialogue systems using degrees of groundedness](#). In *AAAI Conference on Artificial Intelligence*.
- Omar Shaikh, Kristina Gligorić, Ashna Khetan, Matthias Gerstgrasser, Diyi Yang, and Dan Jurafsky. 2024. [Grounding gaps in language model generations](#). *Preprint*, arXiv:2311.09144.
- Aditya K Surikuchi, Sandro Pezzelle, and Raquel Fernández. 2023. [GROOVIST: A metric for grounding objects in visual storytelling](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3331–3339, Singapore. Association for Computational Linguistics.
- Stefanie Tellex, Thomas Kollar, Steven Dickerson, Matthew Walter, Ashis Banerjee, Seth Teller, and Nicholas Roy. 2011. [Understanding natural language commands for robotic navigation and mobile manipulation](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 25(1):1507–1514.
- Yapeng Tian, Di Hu, and Chenliang Xu. 2021. Cyclic co-learning of sounding object visual grounding and sound separation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2745–2754.
- Deepak Venugopal and Vibhav Gogate. 2014. [Scaling-up importance sampling for markov logic networks](#). In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.
- Ning Wang, Jiajun Deng, and Mingbo Jia. 2023a. [Cycle-consistency learning for captioning and grounding](#). *Preprint*, arXiv:2312.15162.
- Renhao Wang, Jiayuan Mao, Joy Hsu, Hang Zhao, Jiajun Wu, and Yang Gao. 2023b. [Programmatically grounded, compositionally generalizable robotic manipulation](#). *Preprint*, arXiv:2304.13826.
- Zehan Wang, Haifeng Huang, Yang Zhao, Linjun Li, Xize Cheng, Yichen Zhu, Aoxiong Yin, and Zhou Zhao. 2023c. [3drp-net: 3d relative position-aware network for 3d visual grounding](#). *Preprint*, arXiv:2307.13363.
- Syed Talal Wasim, Muzammal Naseer, Salman Khan, Ming-Hsuan Yang, and Fahad Shahbaz Khan. 2024. Videogrounding-dino: Towards open-vocabulary spatio-temporal video grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18909–18918.
- Linhui Xiao, Xiaoshan Yang, Xiangyuan Lan, Yaowei Wang, and Changsheng Xu. 2024. [Towards visual grounding: A survey](#). *Preprint*, arXiv:2412.20206.
- Li Yang, chunfeng yuan, Ziqi Zhang, Zhongang Qi, Yan Xu, Wei Liu, Ying Shan, Bing Li, Weiping Yang, Peng Li, Yan Wang, and Weiming Hu. 2023. [Exploiting contextual objects and relations for 3d visual grounding](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 49542–49554. Curran Associates, Inc.
- Zhengyuan Yang, Songyang Zhang, Liwei Wang, and Jiebo Luo. 2021. Sat: 2d semantics assisted training for 3d visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1856–1866.
- Yunan Zeng, Yan Huang, Jinjin Zhang, Zequn Jie, Zhenhua Chai, and Liang Wang. 2024. Investigating compositional challenges in vision-language models for visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14141–14151.
- Hao Zhang, Hongyang Li, Feng Li, Tianhe Ren, Xueyan Zou, Shilong Liu, Shijia Huang, Jianfeng Gao, Lei Zhang, Chunyuan Li, and Jianwei Yang. 2023. [Llava-grounding: Grounded visual chat with large multimodal models](#). *Preprint*, arXiv:2312.02949.
- Zhu Zhang, Zhou Zhao, Zhijie Lin, jieming zhu, and Xiquang He. 2020. [Counterfactual contrastive learning for weakly-supervised vision-language grounding](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 18123–18134. Curran Associates, Inc.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. [Gpt-4v\(ision\) is a generalist web agent, if grounded](#). *Preprint*, arXiv:2401.01614.
- Li Zhou, Zikun Zhou, Kaige Mao, and Zhenyu He. 2023. Joint visual grounding and tracking with natural language specification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 23151–23160.
- Chenming Zhu, Tai Wang, Wenwei Zhang, Kai Chen, and Xihui Liu. 2024. [Scanreason: Empowering 3d visual grounding with reasoning capabilities](#). *Preprint*, arXiv:2407.01525.

A Distribution of years per conference

In Fig 2 and Fig 3, we only showed what the aggregated counts and frequency of “grounding” over time were, possibly hiding some trends. The following figures showcase that the individual conference trends follow the overall trend of increasing in both count and frequency over time.

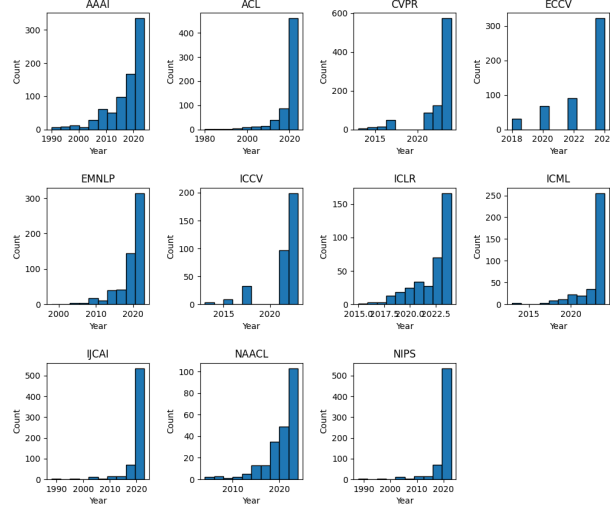


Figure 20: Count of “grounding” for all conferences over time

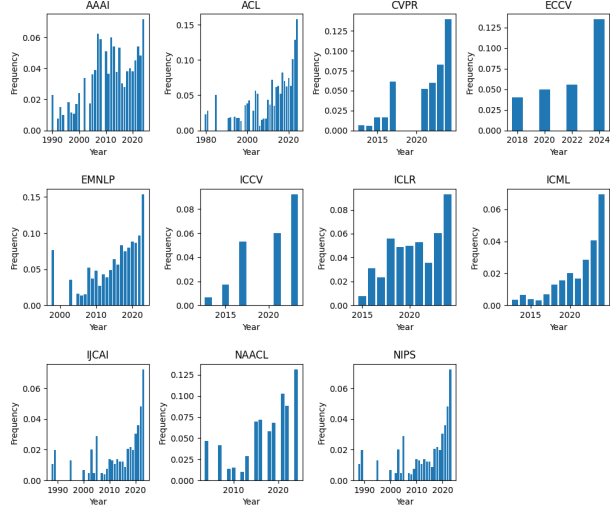


Figure 21: Frequency of papers containing the word “grounding” for all conferences over time

B Words to match per word sense

For counting co-occurrences in section 4.2, we want to pattern match for specific words. In particular, we will match a single word from the list while doing partial matching, and being case insensitive. We ensure that each paper is only counted once. The following table shows which words we pattern match to for each different word sense. For certain words such as “action”, we decided to add on “grounding” to ensure that the commonly used word is not simply just a misinterpretation of the word count. Similarly, we also ensured that the count for each reference word should be at least 3 or more to reduce the noise.

Word Sense	Words to match
Visual Grounding	visual grounding, image grounding, phrase grounding, referring expression
Action Grounding	action grounding, web, agent, robot
Audio Grounding	audio, sound
Video Grounding	video, spatial, spatio-temporal, temporal, object tracking
3D Grounding	3d, point cloud
Dialogue Grounding	dialogue
Markov Logic Networks Grounding	markov, markov logic networks
Physical Dynamics Grounding	physics, dynamics

Table 2: Counts of unique papers with “grounding” by conference found in the corpora.

C Distribution of word senses per year per conference

The following sections provide the graphs and interpretations for each word sense over the years split by conference.

C.1 Visual Grounding

From Fig 22 and Fig 23, one can infer that the visual grounding problems are related more so to the CVPR and ECCV conferences. This is expected as those conferences deal with computer vision.

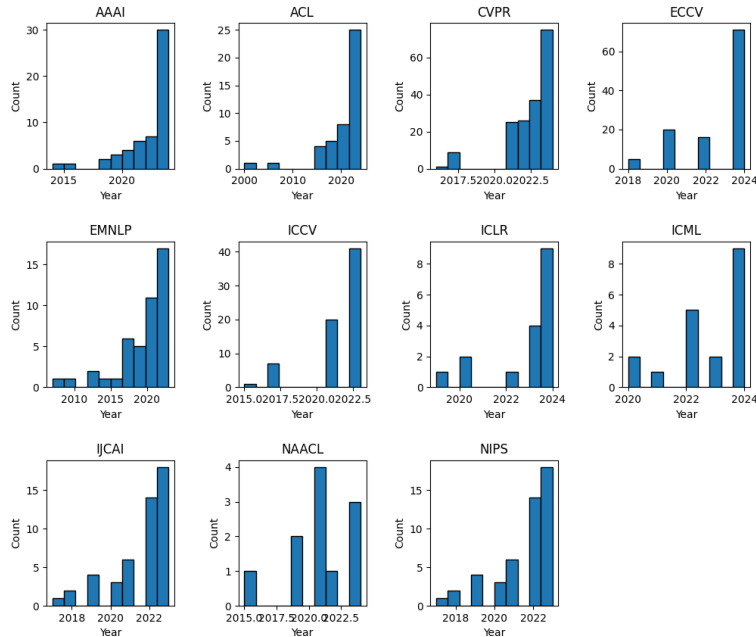


Figure 22: Count of “Visual grounding” per conference over time

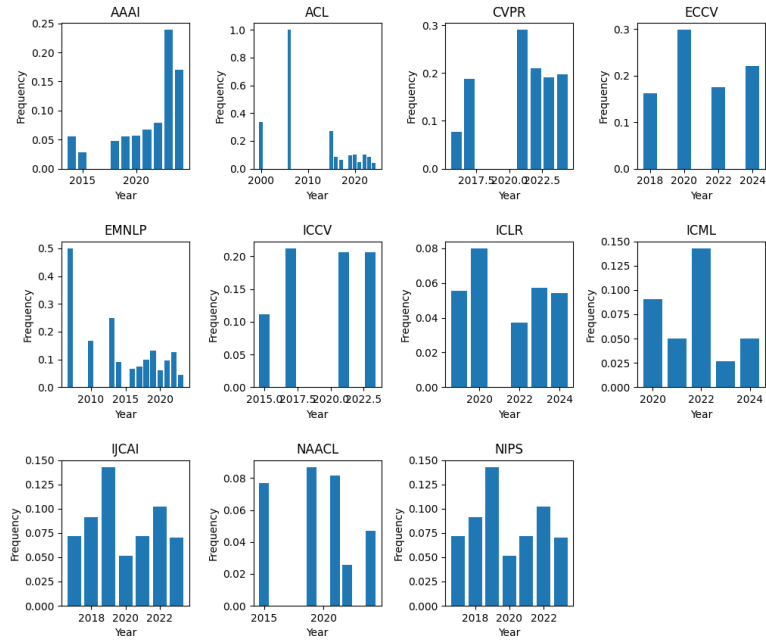


Figure 23: Frequency of "Visual grounding" per conference over time

C.2 Action Grounding

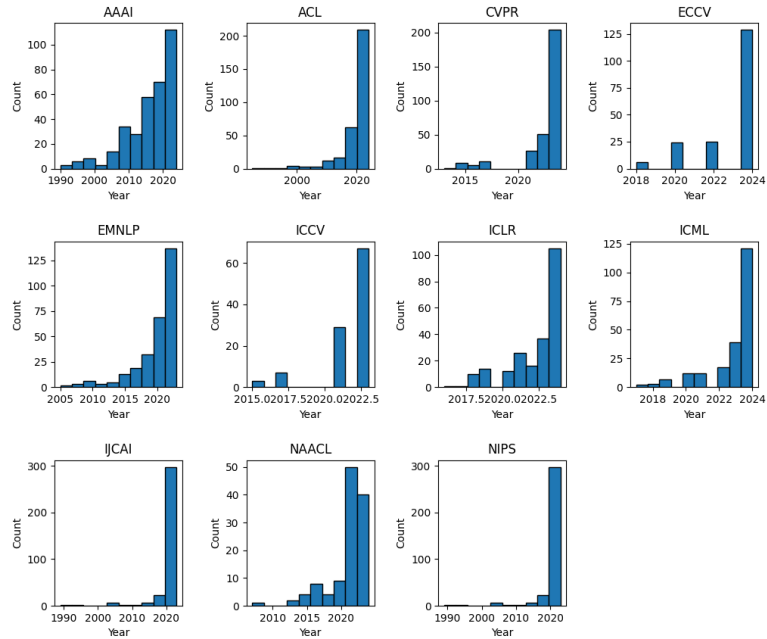


Figure 24: Count of "Action grounding" per conference over time

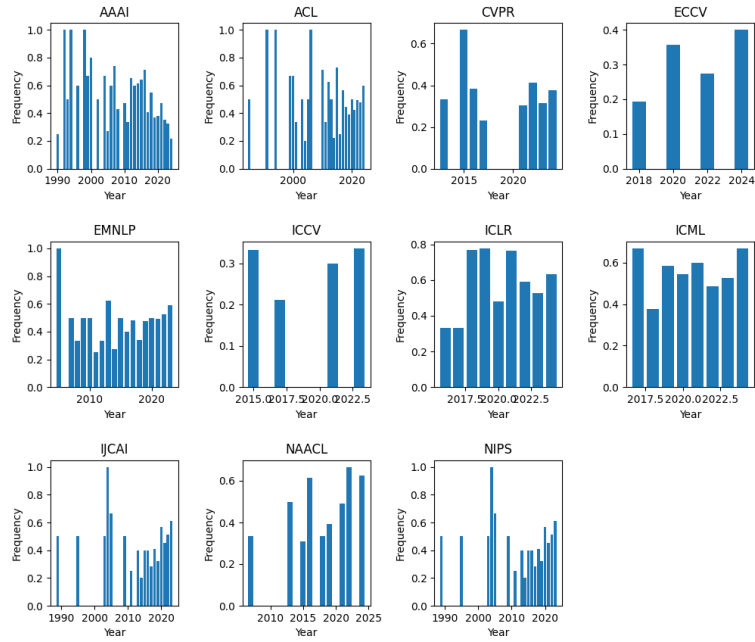


Figure 25: Frequency of “Action grounding” per conference over time

C.3 Audio Grounding

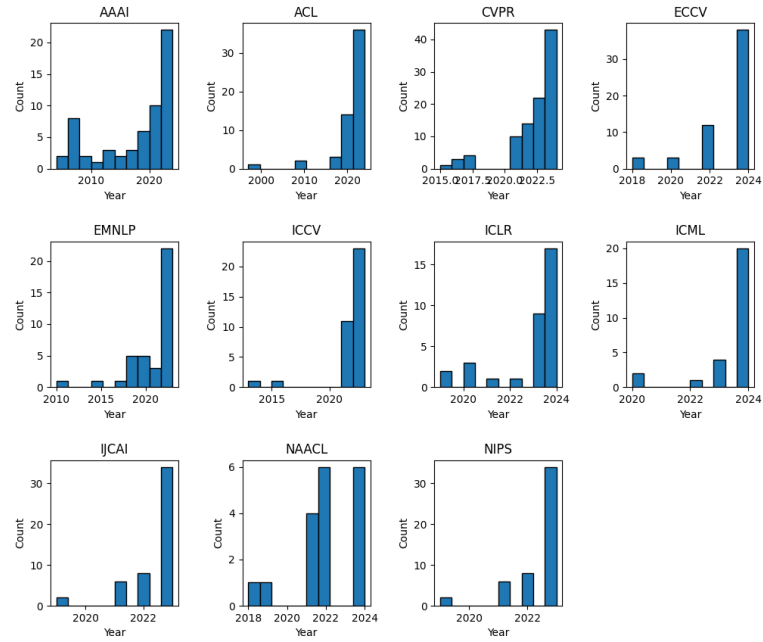


Figure 26: Count of “Audio grounding” per conference over time

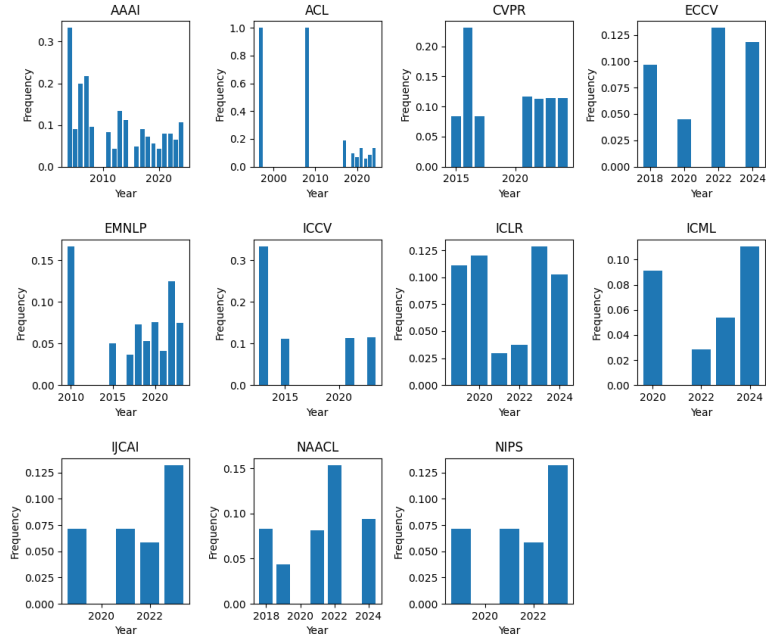


Figure 27: Frequency of “Audio grounding” per conference over time

C.4 Video Grounding

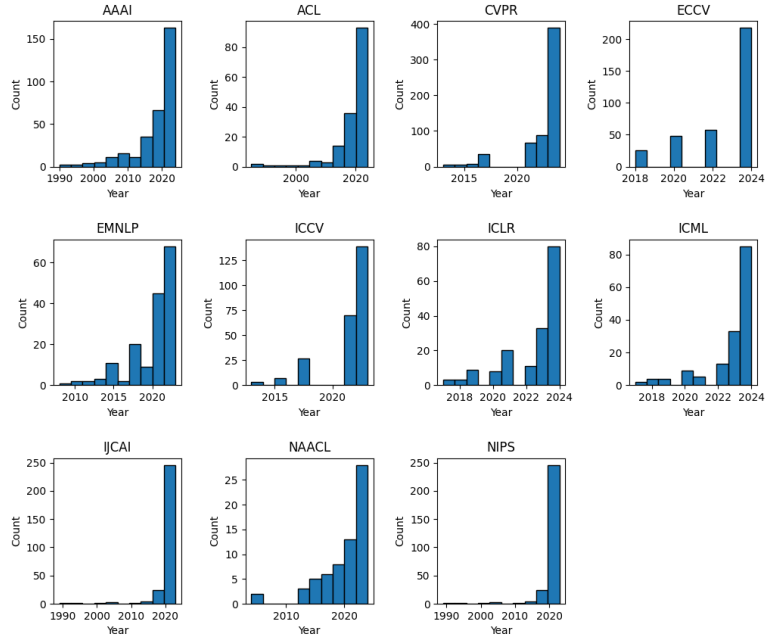


Figure 28: Count of “Video grounding” per conference over time

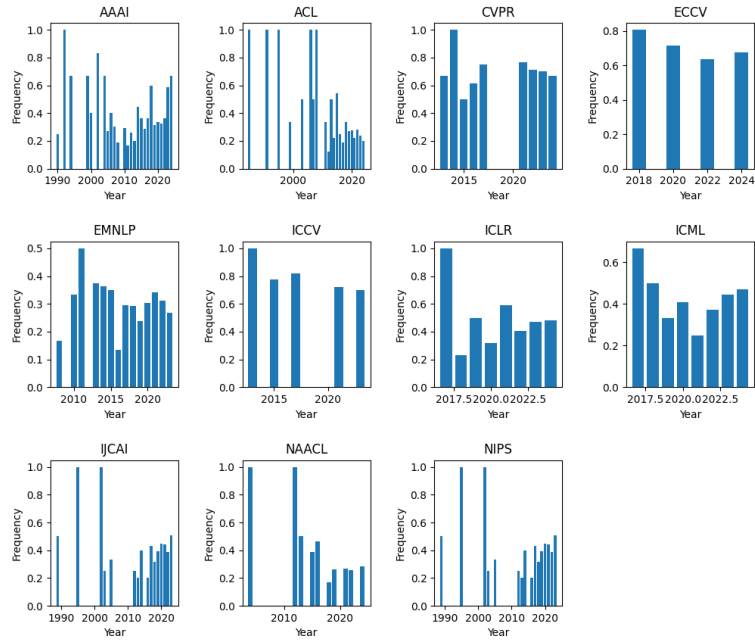


Figure 29: Frequency of “Video grounding” per conference over time

C.5 3D Grounding

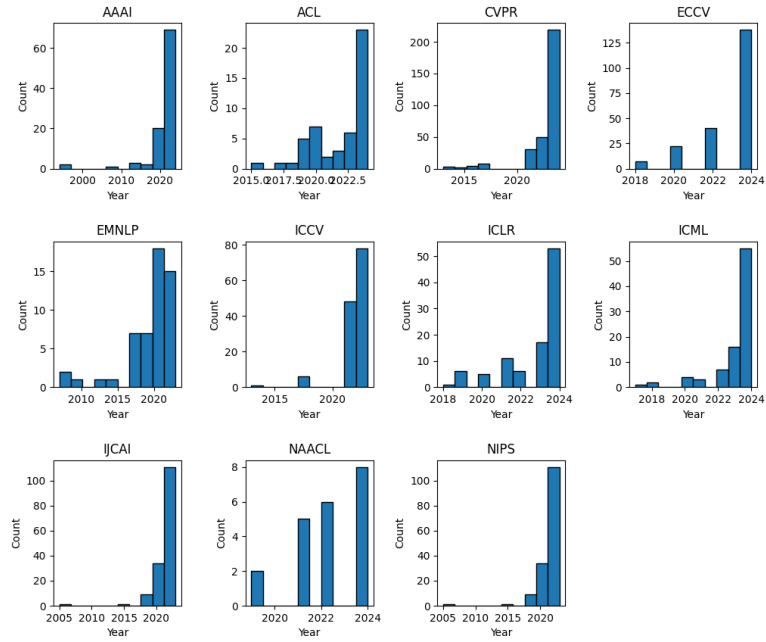


Figure 30: Count of “3D grounding” per conference over time

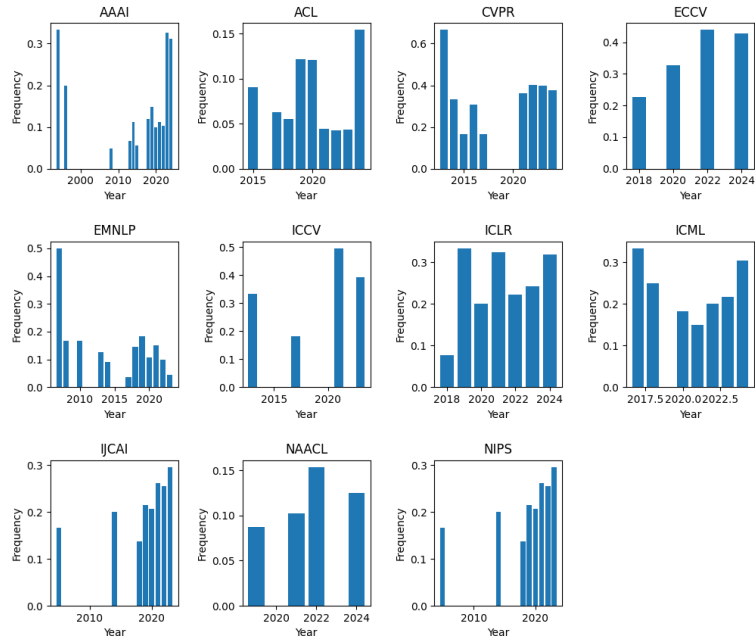


Figure 31: Frequency of “3D grounding” per conference over time

C.6 Dialogue Grounding

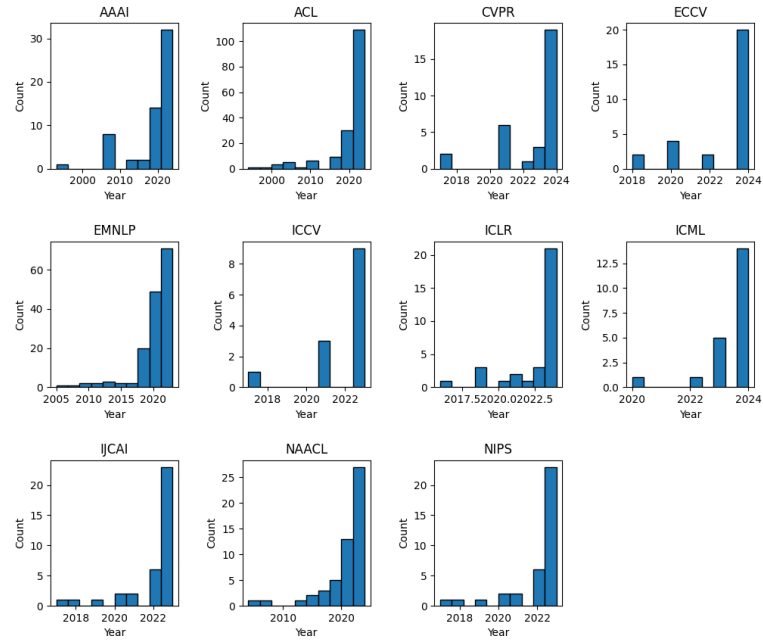


Figure 32: Count of “Dialogue grounding” per conference over time

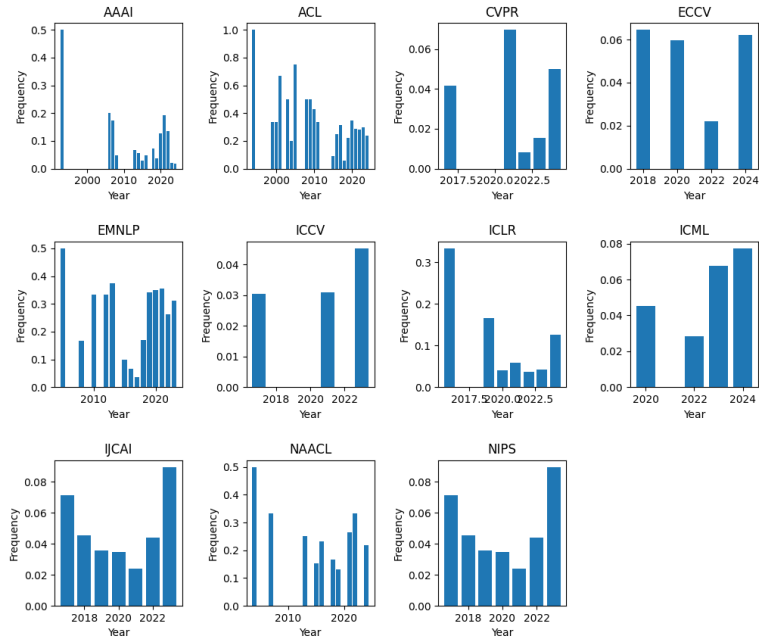


Figure 33: Frequency of “Dialogue grounding” per conference over time

C.7 Markov Logic Networks Grounding

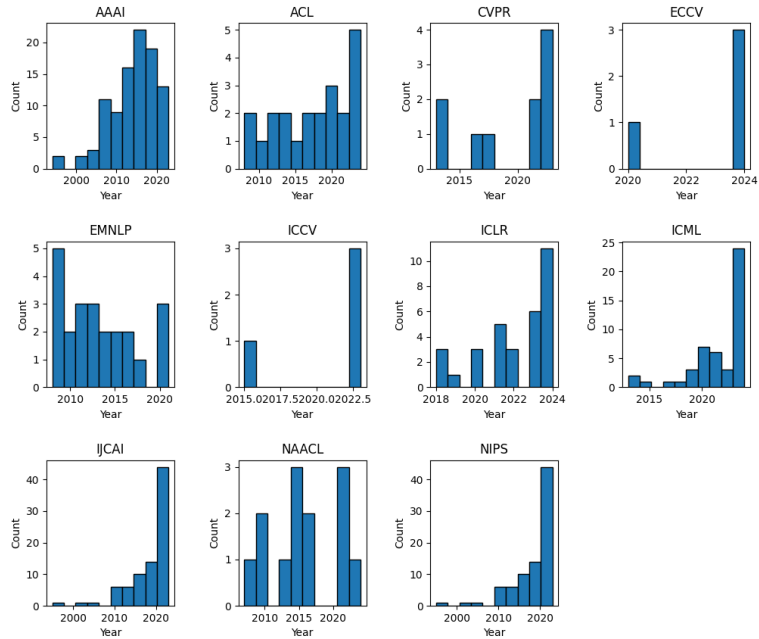


Figure 34: Count of “Markov Logic Networks grounding” per conference over time

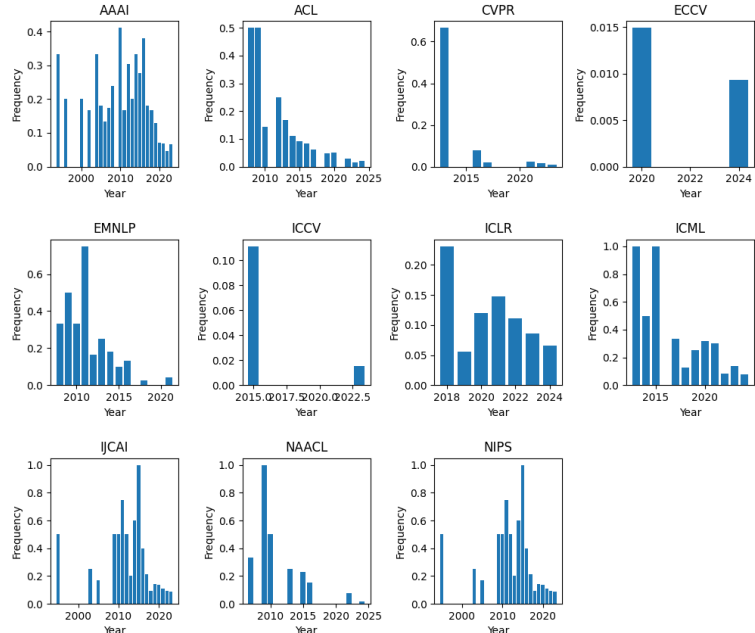


Figure 35: Frequency of “Markov Logic Networks grounding” per conference over time

C.8 Physical Dynamics Grounding

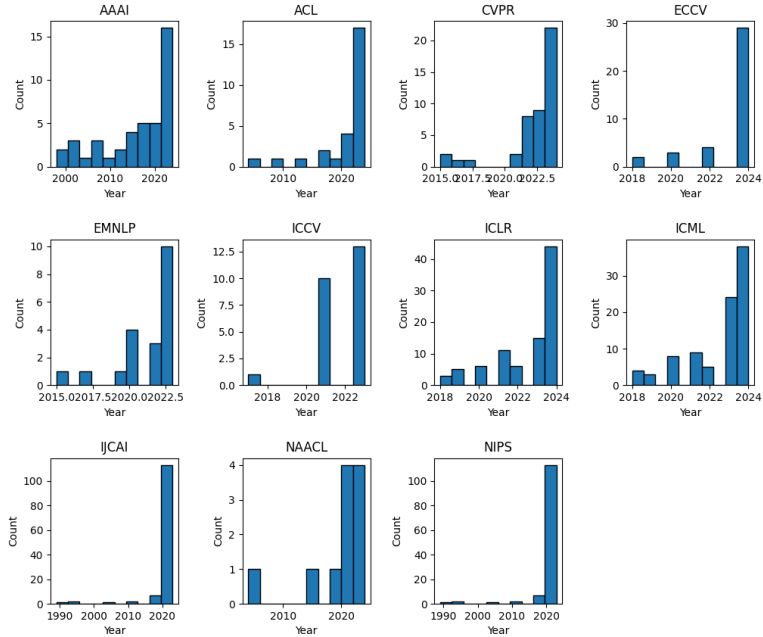


Figure 36: Count of “Physical Dynamics grounding” per conference over time

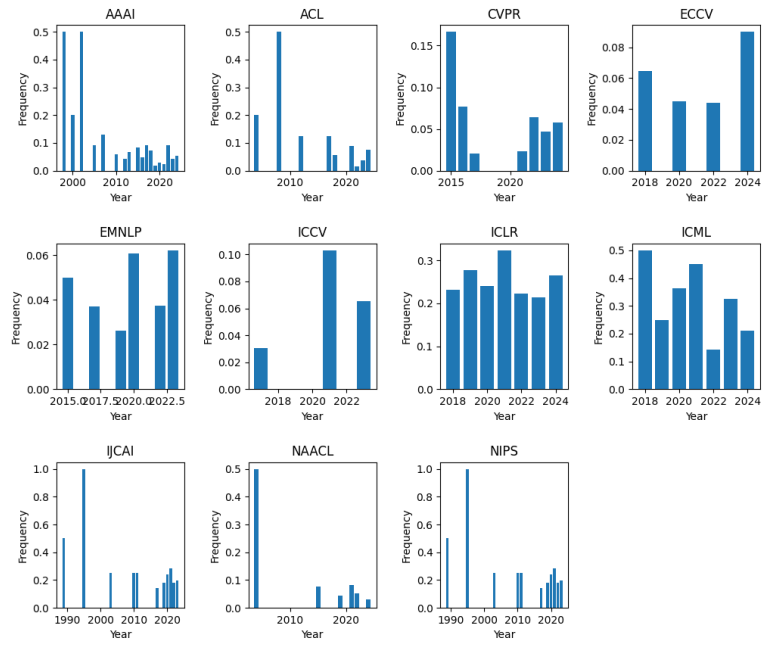


Figure 37: Frequency of "Physical Dynamics grounding" per conference over time