

Language Tasks and Language Games: On Methodology in Current Natural Language Processing Research

David Schlangen

Computational Linguistics / Department of Linguistics

University of Potsdam, Germany

david.schlangen@uni-potsdam.de

Abstract

“This paper introduces a new task and a new dataset”, “we improve the state of the art in X by Y” – it is rare to find a current natural language processing paper (or AI paper more generally) that does *not* contain such statements. What is mostly left implicit, however, is the assumption that this necessarily constitutes progress, and what it constitutes progress towards. Here, we make more precise the normally impressionistically used notions of *language task* and *language game* and ask how a research programme built on these might make progress towards the goal of modelling general language competence.

1 Introduction

Recently, seemingly ever other natural language processing paper introduces a new task and a new dataset.¹ We join some other recent papers (e.g. [Yogatama et al., 2019](#)) in asking whether there is any coherence to this research approach, under which conditions it can lead to progress, and towards what. What we do differently, however, is to look at the fundamental assumptions behind this approach. We try to define central notions, in order to be able to discuss the structure of the typically only implicit formulated approach more clearly.

In our argumentation, we distinguish between *language tasks*, such as for example “describe this image”, or “translate this sentence”—that is, single-step tasks that involve in an essential way natural language material, but not necessarily *only* language material—; *micro worlds*, which are environments that produce disinterested responses to actions, thereby possibly simulating the behaviour of independently existing systems; and *dialogue*

¹Not quite, but not very far. Looking at the 2018 long and short paper proceedings of ACL and EMNLP, we get 94 hits for “introduce new dataset”, 20 hits for “introduce new corpus”, and 101 hits for “introduce new task”.

games as repeated and connected language tasks, which these environments enable. We define these notions first and think about general ways of evaluating their relevance. We close with some tentative recommendations for how to connect individual modelling contributions with the larger enterprise of modelling language processing.

2 Tasks, Worlds, and Games

2.1 Tasks

A *language task* is a mapping between an *input space* and an *output* or *action space*, at least one of which contains natural language expressions.² The mapping has to conform to a *task description*, which is typically given only informally, making reference to theoretical or pre-theoretical constructs external to the definition, such as “translation” or “is true of”. We call this an *intensional description*. Often, a task is also specified *extensionally* through the provision of a *dataset* of examples of the mapping (that is, pairs of state and action), $\mathcal{X} = \{(x_1, y_1), \dots, (x_m, y_m)\}$, where the assumption is that $(x, y) \in \mathcal{X} \rightarrow y = \mathcal{L}(x)$ (\mathcal{L} being the task mapping).

This very general definition essentially covers much, if not all of natural language processing. For example, translation can be seen as a language task where the state space consists of expressions in one language, the action space of expressions in another language, and the task description is that in each pair in the mapping, the second element be a translation of the first. We can further distinguish *understanding tasks*, where the mapping requires demonstration of sensitivity to language meaning (however that is to be further defined); *interpretation tasks*, where the input space contains language expressions that are to be “understood”;

²A formal definition of this and the other notions is given in the Appendix.

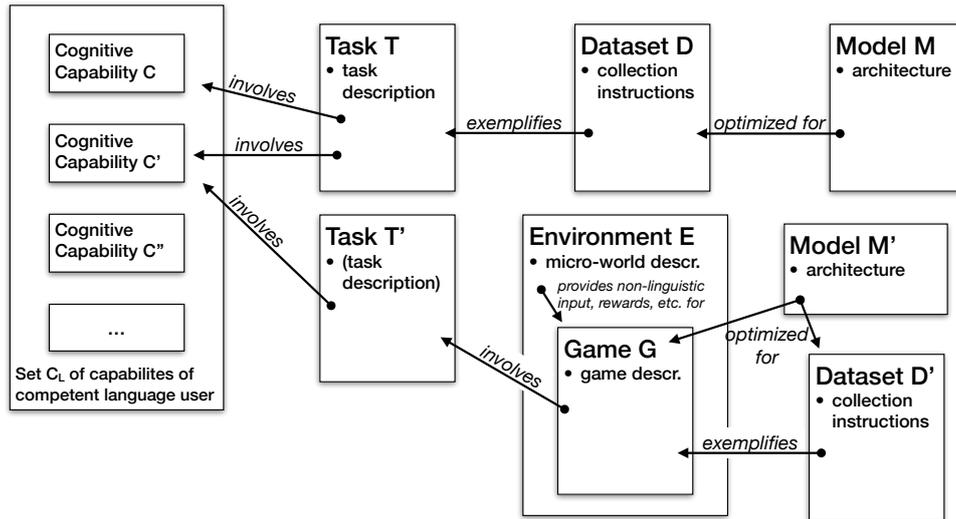


Figure 1: The structure of relations between the research objects *model*, *dataset*, *task*, *game*, *environment*, *cognitive capability*.

generation tasks, where the output does (with a given task potentially being both an interpretation and a generation task); *reference tasks*, where the understanding is shown by relating linguistic and non-linguistic material, and *inference tasks*, where linguistic material is related.

2.2 Worlds

The language competence of humans plays out in *repeated* task, not single-step ones as described in the previous section, and it plays out in contexts where language use is embedded in a non-linguistic context. To study such repeated, situated games, much recent work has made use of environment simulators that compute reactions to actions performed within them, in accordance to the assumed (or actual) rules of the domain they represent.³ Such environments can again be described as mappings, in this case from an action to an environmental response (a state), where again the mapping conforms to a description of which real-world counterpart it is intended to model, and how.⁴

³See for example (Savva et al., 2019; Adams et al., 2012; Johnson et al., 2016; Urbanek et al., 2019; Baroni et al., 2017a; Xia et al., 2018; Yan et al., 2018; Misra et al., 2018; Côté et al., 2018; Bennett and Shatkhin, 2018; Anderson et al., 2018; Savva et al., 2017; Gordon et al., 2017; Brodeur et al., 2017; Chang et al., 2017; Janarthanam and Lemon, 2011; Baroni et al., 2017b; Byron et al., 2007; Yamauchi et al., 2013).

⁴In this desire to model the relevant aspects of a domain, and in the assumption that from dealing with a simulated environment transferable knowledge about dealing with the original environment can be achieved, this approach is reminiscent of the AI microworlds of the 1970s—“we see solving a problem often as getting to know one’s way around a

2.3 Games

An *interaction game* is a setting where *players* can produce *actions* (a special kind of which can be *messages*), possibly regulated by some regime on when they can do this and who can observe them. A privileged, but disinterested player *Nature* can respond to those actions, by providing game-relevant information (and interfacing with the environment in which the game is embedded). Again, we assume that there is an informal *game description* which specifies which, if any, otherwise existing activity the game is meant to model. A *language game* is an interaction game that embeds language tasks which govern the actions of the players. For example, one player asking and the other player answering questions would be a language game that poses repeated language tasks to the players.

Summarising this section, Figure 1 shows the structure of relations between the notions introduced here: Models relate to Tasks, via Datasets, Games are realised in Environments. *Capabilities* will be discussed below.

3 What makes a good task, world, game?

Let’s now assume we encounter a paper that proposes a new dataset, language task, microworld, or language game. How can we evaluate the contribution that is made?

‘micro-world’ in which the problem exists.” (Minsky and Papert, 1972)—and perhaps susceptible to similar kinds of critiques as these attempts (Dreyfus, 1981; Marr, 1982).

3.1 Tasks

... and datasets As mentioned, tasks are often exemplified by the provision of a dataset of examples of the task being executed by agents that are assumed to be capable of doing so—typically, human participants in experiments or data collection efforts. Evaluating such datasets in itself is relatively straightforward. First, it should be *verified*, which is to check whether the provided input/output pairs can indeed be judged correct relative to the task (in its intensional description). If the examples are collected specifically for the purpose of exemplifying the task, this is the process of controlling annotation, and standard methodologies exist (Artstein and Poesio, 2008).

Validating a dataset is a less formalised process. It comprises arguing that the dataset indeed exemplifies the task intension well. For example, pairs only of images of giraffes and sentences describing them would arguably not exemplify the general task of *image description* very well (even if the descriptions are accurate), while perhaps exemplifying the task of *giraffe image description*.

Another way to evaluate datasets is by providing a model of the task learned on parts of it, and testing it on the remaining part (for which a comparison, or *loss*, function on input/output pairs must be provided as well). If a model can “solve” the dataset even when not given information that for theoretical or pre-theoretical reasons is seen to be crucial, the dataset can be considered an unsatisfactory exemplification of the task. E.g., in a *visual (polar) question answering* setting (Antol et al., 2015), if in a dataset all and only the expressions that mention giraffes are true, a model would not need to take the images into account at all to perform well (as it would just need to detect the presence of giraffe-related words), which would be evidence that the dataset is deficient relative to the task description.

... in themselves How can a task in itself be motivated and evaluated? This is easy, if it has a direct value to a consumer (such as translation presumably has), which can be measured. If the consumer is a computer system that processes the output of the task further, the burden of evaluation is simply shifted to the system as a whole. If the interest is in replicating with a theoretically motivated model performance characteristics of humans attempting the task, the task can be evaluated for its power helping distinguish between different mod-

elling choices.

A recent trend, however, has been to motivate tasks in a different way, neither via their inherent practical use, nor as answering questions about language processing as implemented in humans. The argument roughly goes as follows (even if typically only made implicitly): To be good at task T , an agent must possess a set C_T of capabilities (of representational or computational nature). If the $c \in C_T$ are capabilities that competent language users can be shown or argued to possess and make use of in using language—let’s call the set of these capabilities of a competent language user C_L , so that $C_T \subseteq C_L$ —then being able to model these capabilities (via modelling the task) results in progress towards the ultimate goal, which is to model competent language use. And hence, any task T that comes with an interesting set C_T is a good task.⁵

Under what conditions does this argument work? First of all, the assumed connection to the set of capability must indeed be there. We have already seen a way to *challenge* a claimed connection, namely through providing a model that can “solve” a given task (via a dataset) while not having access to information that should be involved in the capability. (Although this challenge

⁵To give some examples of informal versions of this argument, and choosing papers more or less randomly, here are some quotes (typically from the introduction sections of their respective papers):

From the paper that introduced the *visual question answering* task (Antol et al., 2015): “What makes for a compelling AI-complete task? We believe that in order to spawn the next generation of AI algorithms, an ideal task should (i) require multi-modal knowledge beyond a single sub-domain (such as CV) and (ii) have a well-defined quantitative evaluation metric to track progress. [...] Open-ended questions require a potentially vast set of AI capabilities to answer – fine-grained recognition (e.g., What kind of cheese is on the pizza?), object detection (e.g., How many bikes are there?), activity recognition (e.g., Is this man crying?), knowledge based reasoning (e.g., Is this a vegetarian pizza?), and commonsense reasoning (e.g., Does this person have 20/20 vision?, Is this person expecting company?).”

About the *natural language inference* problem, and attempting at least an implicit structuring of the space of capabilities, Condoravdi et al. (2003) write: “The ability to recognize such semantic relations is clearly not a *sufficient* criterion for language understanding: there is more to language understanding than just being able to tell that one sentence follows from another. But we would argue that it is a minimal, *necessary* criterion.”

Williams et al. (2018), on the modern version of this task: “The task of natural language inference (NLI) is well positioned to serve as a benchmark task for research on NLU. [...] In particular, a model must handle phenomena like lexical entailment, quantification, coreference, tense, belief, modality, and lexical and syntactic ambiguity.”

in the first instance only targets the dataset and not the task itself.) Secondly, following usual scientific methodology (Popper, 1934), we can rank the worth of an instantiation of this argument by how precisely the capability is specified, from the trivially correct “task T involves the capability to do task T ” to a statement that could be wrong (and hence involves other theoretical constructs), e.g. “task T involves the capability to compute the syntactic structure of a natural language sentence”.

Furthermore, we can rank the motivation given for a task by how explicit it is in delineating the set of capabilities it involves, along two dimensions. In the one dimension (*separability*), in the most extreme form, the claim would be that the set of capabilities C_T is fully separated from $C_L \setminus C_T$, and hence there is no danger of overfitting solutions to C_T in such a way as would be detrimental for the remaining capabilities.⁶ In the other dimension (*exhaustivity*), the strongest claim would be that T brings out all there is to $c \in C_T$, and that another task T' , insofar as it requires c as well, could be handled by a model of c built with only T in mind. In the other extreme, we only have “ c as required by T ”, which does less to indicate progress beyond T .

As this discussion suggests, it seems difficult to properly motivate a task without relying, at least implicitly, on assumptions about how C_L decomposes.

3.2 Worlds

We have introduced “worlds” (or environments, or simulators) above as the settings that enable repeated tasks (*games*), and in parts their evaluation is connected to those. However, as the providers of *inputs* to a task in reaction to its *outputs*, where this mapping must also confirm to a description, we can also evaluate them qua environment.

First, the notions of *validation* and *verification* apply as well: An environment should match, as well as possible, its stated description of its relation to a real-world counterpart, and should be verified to be correct according to its specification. Environments like those needed for games like Chess and Go (e.g. Silver et al., 2017), or Settlers of Catan (e.g. Afantenos et al., 2012), can fully model their intended real-world counterpart, whereas others can only be modelled approximately

⁶That is, unless the claim is that $C_T = C_L$; the use of “AI-complete” in the quote above from Antol et al. (2015) suggest that is something that not everyone shies away from.

(e.g., “the interior of a house, through which an agent moves, from the perspective of that agent” — as the body of the agent and the agent’s awareness of it, arguably, is a part of the environment, this would entail modelling this as well).

A detailed list of desiderata for such “artificial general intelligence” environments is given by Adams et al. (2012) (see also Baroni et al. (2017b)): “C1. The environment is complex, with diverse, interacting and richly structured objects. C2. The environment is dynamic and open. C3. Task-relevant regularities exist at multiple time scales. C4. Other agents impact performance. C5. Tasks can be complex, diverse and novel. C6. Interactions between agent, environment and tasks are complex and limited. C7. Computational resources of the agent are limited. C8. Agent existence is long-term and continual.” While this list mixes what we separate as demands on environments and on games set in them, it should be useful to evaluate proposed environments.

In analogy to what we observed for tasks, it seems that trying to make progress through modelling tasks in simulated worlds entails making another separability hypothesis, which assumes that the natural competence of handling the world as a whole is separable into handling various parts of it, which can be “knit” together to form the whole.⁷

3.3 Games

What makes for an interesting language game? First we note that games seem to be less well exemplified by datasets than (non-game) tasks are, as for them the relation between an input and an output is much less constrained, and the output can be a sequence of actions rather than a simple one. To give an example, in the language-navigation task (e.g. Anderson et al., 2018; Ma et al., 2019), while the input is a single datum (a verbal description of a goal location), the output is a sequence of navigation actions.⁸

As the site for repeated and connected language

⁷“[W]e feel [the micro-worlds] are so important that we plan to assign a large portion of our effort to developing a collection of these micro-worlds and finding how to embed their suggestive and predictive powers in larger systems without being misled by their incompatibility with literal truth. [...] In order to study such problems, we would like to have collections of knowledge for several ‘micro-worlds’, ultimately to learn how to knit them together.” (Minsky and Papert, 1972).

⁸On the role of environments/games vis-à-vis datasets, Savva et al. (2019) state their belief that “simulators will assume the role played previously by datasets”; no further argumentation is given to support this belief, however.

tasks, we can evaluate a game again for how it connects to language capabilities, and for how the game setting improves over a (non-repeated) task setting. For example, some previous work has shown that language production under (interactive) task constraints is different from null-context language production (Ilinykh et al., 2018; da Silva Rocha and Paraboni, 2016).

Finally, it seems that for language games there is a natural supremum, which would be “unrestricted situated language interaction”. We can then also evaluate proposed tasks for how close they come to this ultimate form, for some definition of “closeness”.

4 Making Progress

Überhaupt hat der Fortschritt das an sich, daß er viel größer ausschaut als er wirklich ist. (It is in the nature of progress that it appears much greater than it actually is.) NESTROY, via Wittgenstein (1953/84)

Research is an incremental enterprise, and new tasks, datasets, models, environments, and games are introduced in a context of existing ones. We can now distinguish several modes of making progress in the general project, by relating new proposals to existing ones.

4.1 Better Models

This is the mode of progress for most of the current work in NLP / AI, and it is also the one that needs the least argumentative support: If a model trained on a given dataset performs better, in terms of the same pre-defined metrics, than a previous model, then progress has been made. Or would that it were so simple: Even in this setting, considerations of computing power spent should arguably also factor, and how to account for use of larger datasets (for example in pre-training), when judging models that achieve better results.⁹ And ultimately, the worth of progress in modelling a particular task rests only on the worth of the task itself.

A different way in which models can improve over others is in how well they are suited for transfer to other tasks (see Ruder (2019) for a recent overview of such transfer-learning).

⁹See e.g. the recent discussion on <https://hackingsemantics.xyz/2019/leaderboards/>.

4.2 Better Datasets

We have discussed the relation between tasks and datasets above. If a dataset can be shown to be successfully modelled even in the absence of information that is deemed critically involved in the capability of interest, it can be considered invalid relative to the task description, and, if the interest in the task is kept, a new dataset must be found.

The task of visual question answering provides an interesting example case of such a development. After Antol et al. (2015) introduced the first large scale dataset for this task, it quickly became clear that this dataset could be handled competitively by models that were deprived of visual input (“language bias”, as noted e.g. by Jabri et al., 2016). This problem was then addressed by Goyal et al. (2017) with the construction of a less biased (and hence more valid) corpus. Targetting the set of capabilities involved in the task, Andreas et al. (2016) noted that “questions in most existing natural image datasets are quite simple, for the most part requiring that only one or two pieces of information be extracted from an image in order to answer it successfully”, which makes it not challenging enough in terms of “[c]ompositionality, and the corresponding ability to answer questions with arbitrarily complex structure”. To improve on that, they introduced the SHAPES dataset which pairs synthetic images with synthetic, programmatically generated sentences that contain spatial relations. Two more datasets explored this direction (Johnson et al., 2017; Suhr et al., 2017), until another dataset (Suhr et al., 2018) progressed beyond the use of synthetic images, in effect claiming that natural images are more valid for the task described as “visual question answering”.

4.3 Better Tasks and Games

We see a movement towards involvement of *capabilities* already in the previous two sections. Where a model can be judged to be better than another one (trained on the same data) simply by the linear order imposed by the evaluation metric, a dataset must be argued to be more representative of a task by taking recourse to the capabilities of interest (e.g., in the example discussed just above, that of handling *compositionality*).

As shown in Figure 1, tasks are only grounded (to their left in the diagram) by capabilities (and games by tasks), and hence an argumentation for a task T' in relation to a previous task T should

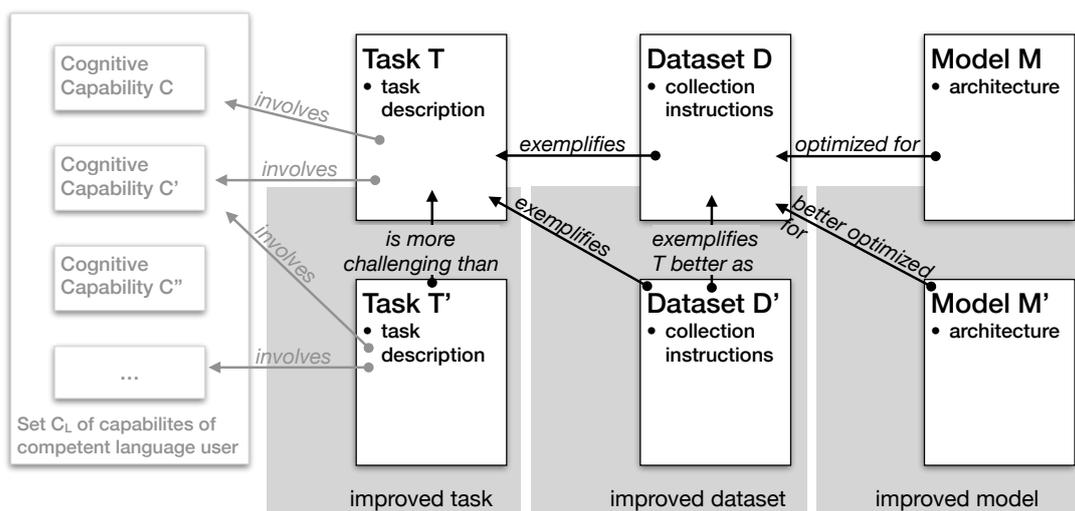


Figure 2: Three ways of making progress: improving models, improving datasets, improving tasks. (Not shown: devising models that handle more than one dataset.)

ideally make mention of how $C_{T'}$ and C_T relate. (Similarly for games.)

Figure 2 again summarises this discussion in a diagram.

4.4 From Models to Capabilities

An interesting additional avenue of research has been explored in recent years (see for example the paper by [Hewitt and Manning, 2019](#), and references therein), where models trained on datasets (representing tasks or games) are explored for how they decomposed their (representational) task. For example, the question might be, as in the cited paper, whether a particular model trained on a particular task creates “internally” something akin to syntax trees. Again, why this might be interesting is typically not spelled out in these papers, but one can assume that the intended underlying argument goes something like this: “If the theoretical construct is to be found, this shows that postulating it is empirically validated (and it can be learned simply through exposure to data); if it is not there, then the task can be done without it, to the extent that the model can do it, and it need not be postulated”.

5 Conclusions: Making it Explicit

The discussion above has mixed normative (how it could, or should be) and descriptive (how it is) aspects. I will close by commenting more directly on what I think could be improved.

As discussed above, the most frequent way of

contributing to the project is by improving models, for which established methods of evaluation exist. Next frequent is providing new or improved datasets. Here already we find a much larger variety of how the contribution is motivated and framed, and often it is taken as self-evident that a new dataset will drive progress. Here, more explicitness about the assumed advantages, and how they connect to the goal of modelling language competence, would be helpful.¹⁰

This holds even more when introducing new tasks or games. As the discussion above hopefully has made plausible, motivating those in themselves, and relating them to existing ones, requires making claims about language capabilities. To make these in a convincing way, it seems advantageous to strive for a renewed, stronger connection to the sciences that study them: linguistics and cognitive psychology. Those fields provide the theoretical terms and constructs that would allow us to make the, as discussed above, typically implicit arguments more explicit (and hence contestable).

Acknowledgements

I thank Raquel Fernández, Manfred Stede, and Sina Zarrieß for interesting discussions about topics related to this paper, while reserving the exclusive right to be blamed for any misunderstandings it might betray.

¹⁰The recent initiative of formulating “Data Sheets” ([Geburu et al., 2018](#)), although developed for different purposes, would be a good step also in this direction.

References

- Sam Adams, Itmar Arel, Joscha Bach, Robert Coop, Rod Furlan, Ben Goertzel, J. Storrs Hall, Alexei Samsonovich, Matthias Scheutz, Matthew Schlesinger, Stuart C. Shapiro, and John Sowa. 2012. [Mapping the Landscape of Human-Level Artificial General Intelligence](#). *AI Magazine*, 33(1):25–42.
- Stefanos Afantenos, Nicholas Asher, Farah Benamara, Anais Cadilhac, Cedric Degremont, Pascal Denis, Markus Guhe, Simon Keizer, Alex Lascarides, Oliver Lemon, Philippe Muller, Soumya Paul, Verena Rieser, and Laure Vieu. 2012. Developing a corpus of strategic conversation in the settlers of catan. In *Proceedings of the 1st Workshop on Games and NLP*, Kanazawa, Japan.
- Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. 2018. [Vision-and-Language Navigation: Interpreting visually-grounded navigation instructions in real environments](#). In *CVPR 2018*.
- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. Neural Module Networks. In *Proceedings of the 2016 Conference on Computer Vision and Pattern Recognition (CVPR 2016)*.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *International Conference on Computer Vision (ICCV)*.
- Ron Artstein and Massimo Poesio. 2008. [Inter-Coder Agreement for Computational Linguistics](#). *Computational Linguistic*, 34(4):555–596.
- Marco Baroni, Armand Joulin, Allan Jabri, Germàn Kruszewski, Angeliki Lazaridou, Klemen Simoncic, and Tomas Mikolov. 2017a. [CommAI: Evaluating the first steps towards a useful general AI](#). *arXiv*, pages 1–9.
- Marco Baroni, Claes Strannegård, David L. Dowe, Katja Hofmann, Kristinn R. Thórisson, Jordi Bieger, Nader Chmait, Fernando Martínez-Plumed, and José Hernández-Orallo. 2017b. [A New AI Evaluation Cosmos: Ready to Play the Game?](#) *AI Magazine*, 38(3):66.
- Andrew Bennett and Max Shatkhin. 2018. Mapping Instructions to Actions in 3D Environments with Visual Goal Prediction. In *EMNLP 2018*, pages 2667–2678.
- Simon Brodeur, Ethan Perez, Ankesh Anand, Florian Golemo, Luca Celotti, Florian Strub, Jean Rouat, Hugo Larochelle, and Aaron Courville. 2017. [HoME: a Household Multimodal Environment](#). *ArXiv*.
- Donna Byron, Alexander Koller, Jon Oberlander, Laura Stoia, and Kristina Striegnitz. 2007. Generating Instructions in Virtual Environments (GIVE): A Challenge and an Evaluation Testbed for NLG. In *Proceedings of the Workshop on Shared Tasks and Comparative Evaluation in Natural Language Generation*.
- Angel Chang, Angela Dai, Thomas Funkhouser, Manolis Savva, and Shuran Song. 2017. [Matterport3D : Learning from RGB-D Data in Indoor Environments](#). *ArXiv*.
- Cleo Condoravdi, Richard Crouch, Valeria de Paiva, Reinhard Stolle, and Daniel G Bobrow. 2003. [Entailment, intensionality and text understanding](#). In *Proc. of the HLT-NAACL 2003 Workshop on Text Meaning*, pages 38–45.
- Marc-Alexandre Côté, Ákos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, Wendy Tay, and Adam Trischler. 2018. [TextWorld: A Learning Environment for Text-based Games](#). *ArXiv*.
- Danillo da Silva Rocha and Ivandr e Paraboni. 2016. Reference production in human-computer interaction : Issues for Corpus-based Referring Expression Generation. In *LREC*, pages 2994–2998.
- Hubert L. Dreyfus. 1981. From micro-worlds to knowledge: AI at an impasse. In John Haugeland, editor, *Mind Design*. MIT Press.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M. Wallach, Hal Daum e III, and Kate Crawford. 2018. [Datasheets for datasets](#). *CoRR*, abs/1803.09010.
- Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. 2017. [IQA: Visual Question Answering in Interactive Environments](#). *ArXiv*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. [Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering](#). In *CVPR 2017*.
- John Hewitt and Christopher D Manning. 2019. [A Structural Probe for Finding Syntax in Word Representations](#). In *NAACL-HLT 2019*.
- Nikolai Ilinykh, Sina Zarri e, and David Schlangen. 2018. The Task Matters. Comparing Image Captioning and Task-Based Dialogical Image Description. In *Proceedings of 11th International Conference on Natural Language Generation (INLG 2018)*.
- Allan Jabri, Armand Joulin, and Laurens van der Maaten. 2016. [Revisiting Visual Question Answering Baselines](#). In *European Conference on Computer Vision (ECCV)*.

- Srini Janarthanam and Oliver Lemon. 2011. [The GRUVE Challenge : Generating Routes under Uncertainty in Virtual Environments](#). In *ENLG '11 Proceedings of the 13th European Workshop on Natural Language Generation*, pages 208–211.
- Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. 2017. [CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning](#). In *CVPR 2017*, pages 1988—1997.
- Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. 2016. [The malmo platform for artificial intelligence experimentation](#). *IJCAI International Joint Conference on Artificial Intelligence*, 2016-Janua:4246–4247.
- Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan Al-Regib, Zsolt Kira, Richard Socher, and Caiming Xiong. 2019. [Self-Monitoring Navigation Agent via Auxiliary Progress Estimation](#). *ArXiv*, pages 1–18.
- David Marr. 1982. *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. W.H. Freeman, San Francisco, USA.
- Marvin Minsky and Seymour Papert. 1972. [Progress Report on Artificial Intelligence](#). Technical report, MIT Artificial Intelligence Laboratory, Cambridge, Mass., USA.
- Dipendra Misra, Andrew Bennett, Valts Blukis, Eyvind Niklasson, Max Shatkhin, and Yoav Artzi. 2018. [Mapping Instructions to Actions in 3D Environments with Visual Goal Prediction](#). *ArXiv*.
- Karl Popper. 1934. *Logik der Forschung*. Mohr Siebeck.
- Sebastian Ruder. 2019. *Neural Transfer Learning for Natural Language Processing*. Ph.D. thesis, National University of Ireland, Galway.
- Manolis Savva, Angel X. Chang, Alexey Dosovitskiy, Thomas Funkhouser, and Vladlen Koltun. 2017. [MINOS: Multimodal Indoor Simulator for Navigation in Complex Environments](#). *ArXiv*, pages 1–14.
- Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. 2019. [Habitat: A Platform for Embodied AI Research](#). *ArXiv*.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George Van Den Driessche, Thore Graepel, and Demis Hassabis. 2017. [Mastering the game of Go without human knowledge](#). *Nature*, 550(7676):354–359.
- Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. 2017. [A Corpus of Natural Language for Visual Reasoning](#). In *Proceedings of the 2017 meeting of the Association for Computational Linguistics (ACL 2017)*.
- Alane Suhr, Stephanie Zhou, Iris Zhang, Huajun Bai, and Yoav Artzi. 2018. [A Corpus for Reasoning About Natural Language Grounded in Photographs](#). In *Proceedings of NIPS 2018*, Montreal, Canada.
- Bernard Suits. 1978. *The Grasshopper: Games, Life, and Utopia*. The University of Toronto Press, Toronto, Canada.
- Richard S. Sutton and Andrew G. Barto. 1998. *Reinforcement Learning*. MIT Press, Cambridge, USA.
- Jack Urbanek, Angela Fan, Siddharth Karamcheti, Saachi Jain, Samuel Humeau, Emily Dinan, Tim Rocktäschel, Douwe Kiela, Arthur Szlam, and Jason Weston. 2019. [Learning to Speak and Act in a Fantasy Text Adventure Game](#). *ArXiv*.
- Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. [A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference](#). In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Ludwig Wittgenstein. 1953/84. *Tractatus Logicus Philosophicus und Philosophische Untersuchungen*, volume 1 of *Werkausgabe*. Suhrkamp, Frankfurt am Main.
- Fei Xia, Amir Zamir, Zhi-Yang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. 2018. [Gibson Env: Real-World Perception for Embodied Agents](#). In *CVPR 2018*.
- Takashi Yamauchi, Mikio Nakano, and Kotaro Funakoshi. 2013. [A Robotic Agent in a Virtual Environment that Performs Situated Incremental Understanding of Navigational Utterances](#). In *SIGdial 2013*, August, pages 369–371.
- Claudia Yan, Dipendra Misra, Andrew Bennett, Aaron Walsman, Yonatan Bisk, and Yoav Artzi. 2018. [CHALET : Cornell House Agent Learning Environment](#). *ArXiv*.
- Dani Yogatama, Cyprien de Masson D’Autume, Jerome Connor, Tomas Kocisky, Mike Chrzanowski, Lingpeng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu, Chris Dyer, and Phil Blunsom. 2019. [Learning and Evaluating General Linguistic Intelligence](#). *ArXiv*, pages 1–14.

A Formalising the Notions

A.1 Tasks

Definition 1 A Language Task is a tuple (S, A, \mathcal{L}, D_T) , where:

- S is a (possibly infinite) set of states,
- A is a (possibly infinite) set of actions,
- with either the states in S or the actions in A (or both) having as part natural language expressions, and
- $\mathcal{L} : S \rightarrow A$ is a function that maps a state $s \in S$ to an action $a \in A$, where
- the mapping \mathcal{L} conforms to task description D_T .

A.2 Worlds

Definition 2 A Micro-World or Environment is a tuple $(S, A, \mathcal{E}, R, D_W)$, where function $\mathcal{E} : S \times A \rightarrow S \times R$ maps an action a , taken in state s , to a state s' and a reward r , and the mapping conforms to the world description D_W .

A.3 Games

We approach the definition of a *language game* via the more general notion of *interaction game*:¹¹

Definition 3 An Interaction Game is a tuple (P, A, o, T, E, D_G) , where:

- $P = \{p_1, \dots, p_n, N\}$ is the set of n regular players p , together with one additional player N (for Nature).
 N has a special status in that it does not have a strategic interest in the outcome of the game.
- $A = \{A_1, \dots, A_n, A_N\}$ is the set of action spaces, with one space per player.
Action types can be complex: e.g., (nav, s) , for “navigation action, south”, or $(utt, \text{“I don’t know”})$ for “utterance of I don’t know”. If defined in the right way (for example by a recursive grammar), the set of action types can be infinite. Players chose actions from their space of available actions; the resulting action tokens a_j are associated with their originator through a

¹¹Pace Wittgenstein (1953/84), we get to define what a game is. Or we could go with Suits (1978), who is happy to define games as rule-guided activities of voluntary attempt to overcome unnecessary obstacles. His concepts of *preludory goal*, which can be stated independently of the game (e.g., in football (soccer), “make the ball be in the opponent team’s goal”); *constitutive rules*, which make reaching that goal more difficult than necessary (e.g., by disallowing to just grab the ball and carry it to the goal); and *lusory attitude*, which is to accept the complications posed by the constitutive rules, can inform the design of games that work via crowdsourcing.

function $a : A_i \rightarrow P$, and with a position in the sequence of actions that have been performed since the beginning of the interaction through a function $t : A_i \rightarrow \mathbb{N}$.

- $o : P_i \times A_i \rightarrow \mathcal{P}(P)$ (for $P_i \in P, A_i \in A$) is the observability function that specifies which types of actions by which player can be observed by which subset of the players.

In normal cases, one would assume that players can observe their own actions, and that Nature observes all actions; but this allows for the specification of deviant cases.

- $T : \emptyset \cup (P \times A) \rightarrow \mathcal{P}(P)$ is the turn taking rule that specifies who can act next, depending on who did what last. It also specifies who can start the game.

In a *free initiative* setting, any player can act at any time; in a strict turn based setting, the current player would always be excluded from the set of next players.

- $E : S \rightarrow V$ is the evaluation rule that maps a sequence of action tokens $\langle a_1, \dots, a_m \rangle$ into an evaluation, where the set of possible evaluations V includes at least one positive one (e.g., success) and one negative one (e.g., failure).

The evaluation is made known to the players when a positive or negative outcome has been reached. If it is not, or if it does not contain outcomes denoted as positive or negative, we call the resulting structure an *interaction setting*, rather than an interaction game.

- D_G finally is the game description which specifies which, if any, otherwise existing activity the game is meant to approximate.

The well-known *Gridworld* game (see e.g. Sutton and Barto (1998)) for example can be represented in these terms as being an interaction game with one regular player (the agent) interacting with Nature, $P = \{p_1, N\}$. The agent can only perform navigation actions: $A_1 = \{(nav, n), \dots\}$, Nature informs on the resulting available navigation options, $A_N = \{(inform, (n, w)), \dots\}$, with the information that Nature relays coming from a microworld that simulates the grid and the movement on it.

Gridworld does not involve language, and hence is not an example of a *language* interaction game. As an example of a language interaction

setting that is not a game, we can define *free chat interaction* in our terms as involving two players and an inert Nature that does not intervene: $P = \{p_1, p_2, N\}$, $A_1 = A_2 = \{(utt, \alpha)\}$, $A_N = \emptyset$, T is a constant function into P (free initiative), all actions are observed by all, E is a constant function into $\{undecided\}$.