



Structuring information in interactive natural language generation

Alexander Koller
koller@ling.uni-potsdam.de

Nikos Engonopoulos
engonopo@uni-potsdam.de

Martín Villalba
mvillalb@uni-potsdam.de

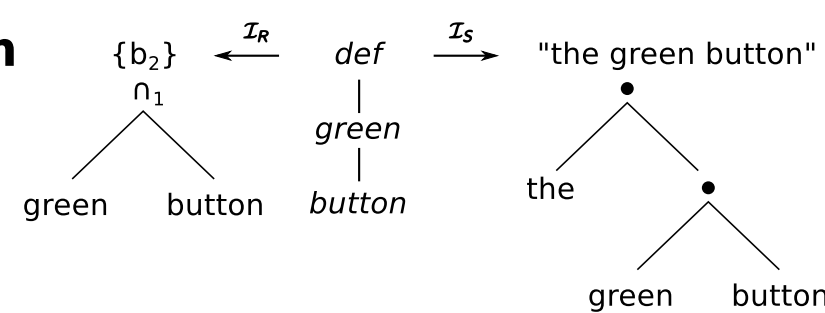
SFB 632 - D6

Research questions

- Given a target object, how to generate the **referring expression** that has the **highest chance** of being correctly understood?
- Given a referring expression and the user's observed behavior, how do we **predict what the user has understood** as its object?
- How do we **recognize that a user has not understood**, and generate the right instruction to **correct it** with appropriate **contrastive focus**?

Generation with semantically interpreted grammars

We use a **synchronous grammar formalism** relating **sets of objects** from a first order model of the world to **natural language strings** describing them, via an **abstract syntactic structure**.



Grammar rule	String	Denotation
NP → def(N)	the • w ₁	uniq(R ₁) = if (R ₁ is singleton) then R ₁ else ∅
N → leftof(N, NP)	w ₁ • to the left of • w ₂	{ a ∈ R ₁ exists b ∈ R ₂ s.t. (a,b) ∈ left_of }
N → green(N)	green • w ₁	[green] ∩ R ₁
N → red(N)	red • w ₁	[red] ∩ R ₁
N → button	button	[button]

We parse the target object into a **finite chart**, a data structure that represents the (possibly infinite) set of all valid referring expressions.

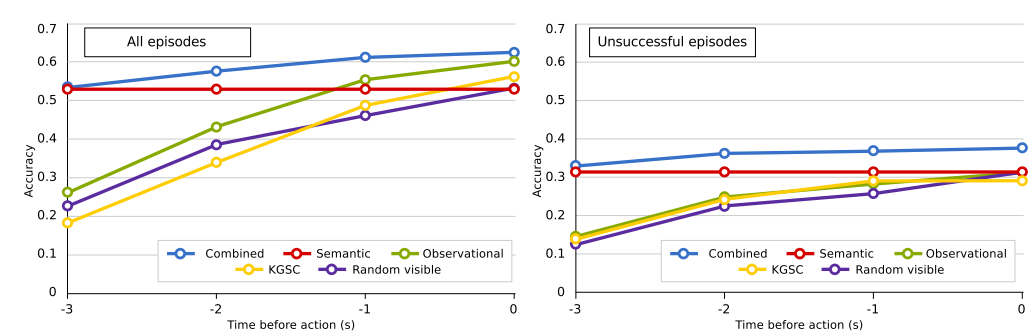
NP/{b ₁ }	→ def(N/{b ₁ })
NP/{b ₂ }	→ def(N/{b ₂ })
N/{b ₁ }	→ leftof(N/{b ₁ , b ₂ , b ₃ }, NP/{b ₂ })
N/{b ₁ , b ₃ }	→ red(N/{b ₁ , b ₂ , b ₃ })
N/{b ₂ }	→ green(N/{b ₁ , b ₂ , b ₃ })
N/{b ₁ , b ₂ , b ₃ }	→ button

We select the **best referring expression** out of the chart, i.e. the one maximizing the probability $P_{sem}(b^*|r, s)$ of b^* being r .

Evaluation

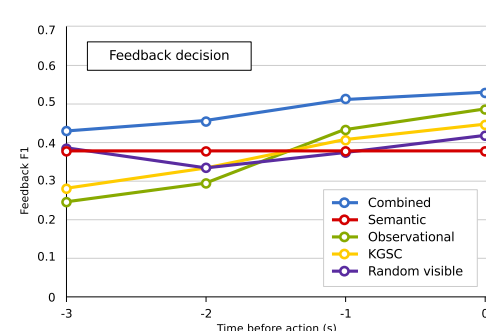
Training data: 3414 episodes for P_{sem} and 6478 tuples $\langle \sigma, b \rangle$ for P_{obs} from the GIVE-2.5 Challenge.
Test data: 5028 episodes from the GIVE-2 Challenge.

The model combining P_{obs} and P_{sem} **outperformed its components and two baselines** on prediction and feedback accuracy.



Prediction accuracy

Predict the object understood by the user: given r, s and σ until some $t > t_0$, predict the object understood by the user: $\text{argmax}_a P(a|r, s, \sigma)$



Feedback accuracy

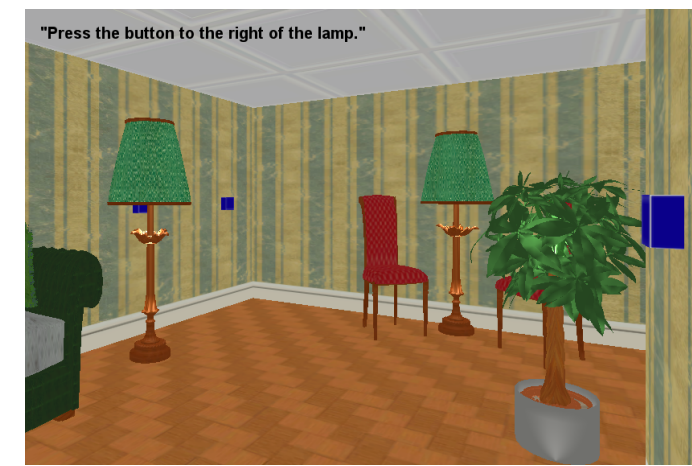
Provide feedback if the system predicts the user is about to make a mistake: given r, s and σ until some $t > t_0$, give feedback if $p(a') - p(a^*) > \theta$ for some object $a' \neq a^*$ ($\theta = 0.1$)

Eye-tracking information has shown to improve accuracy for the P_{obs} model in **harder** episodes. Evaluation of the success and perceived quality of the generation system output with live users using **crowdsourcing** is underway.

Conclusions and ongoing work

- We have developed a **probabilistic semantic model** $P_{sem}(b|r, s)$ that predicts to which object a will the user resolve the RE r in the scene s , and a **probabilistic observational model** $P_{obs}(b|\sigma)$ that updates in real time, based on the user's behavior σ .
- We also presented an algorithm that computes the RE r that has the highest probability of being understood as the target object a .

The GIVE Challenge



Users have to solve a puzzle in a 3D world. They can interact with objects in the world (e.g. click on buttons) and move freely in the environment.

NLG systems guide the users by generating instructions, including referring expressions for objects in the environment.

Grounding problem: Systems have to **predict** (mis)understanding of an object and **prevent** mistakes by providing corrective **feedback**.

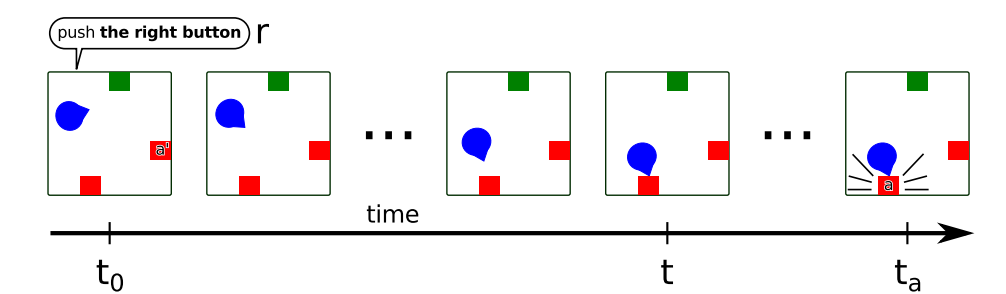
The interaction corpora comprises over 2500 games and more than 340h of recorded interactions.

Probabilistic semantic model

Purpose: predict the probability distribution $P(b|r)$ of the RE r being understood as an object b over the set of all objects in the world.

The **probabilistic semantic model** (P_{sem}) uses information available at *instruction giving time*: syntactic and semantic features of the utterance, potential sources of confusion, visual salience, spatial features (distance, angle, etc.).

Model parameters are automatically **learnt from recorded interactions** using machine learning methods.



An **episode** consists of the recorded events between an instruction and the user's action.

Probabilistic observational model

Purpose: predict the probability distribution $P(b|\sigma)$ of resolving the RE to the object b based on the behavior σ .

The **probabilistic observational model** (P_{obs}) collects new information continuously: eye-tracking, spatial features, visual salience, and the evolution of this information over time. The model **does not rely** on the text of the given instruction.

Like the P_{sem} model, parameters are automatically learnt from recorded interactions.

Combined probabilistic model

When receiving an instruction containing a **referring expression** r at a given **world state** s , the user resolves r to an object b . The user then moves towards b , exhibiting **behavior** σ .

$$P(b | r, s, \sigma) \propto P_{sem}(b|r,s) \cdot P_{obs}(b|\sigma)$$

P_{sem} and P_{obs} are log-linear models, trained separately on recorded interactions and then combined. Both can generalize to unseen worlds.

Further reading:

Nikos Engonopoulos, Martín Villalba, Ivan Titov and Alexander Koller: Predicting the resolution of referring expressions from user behavior. Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP). Seattle, WA, USA. October 2013.

Nikos Engonopoulos and Alexander Koller: Generating effective referring expressions using charts. Proceedings of the 8th International Natural Language Generation Conference (INLG). Philadelphia, PA, USA. June 2014.

Alexander Koller, Kristina Striegnitz, Donna Byron, Justine Cassell, Robert Dale, Johanna Moore and Jon Oberlander: The First Challenge on Generating Instructions in Virtual Environments. In E. Kraemer and M. Theune (eds.), Empirical Methods in Natural Language Generation, volume 5790 of LNAI. Springer, 2010.

Nikolina Koleva, Martín Villalba, Maria Staudte and Alexander Koller: The impact of listener gaze on predicting reference resolution. Currently under review.