

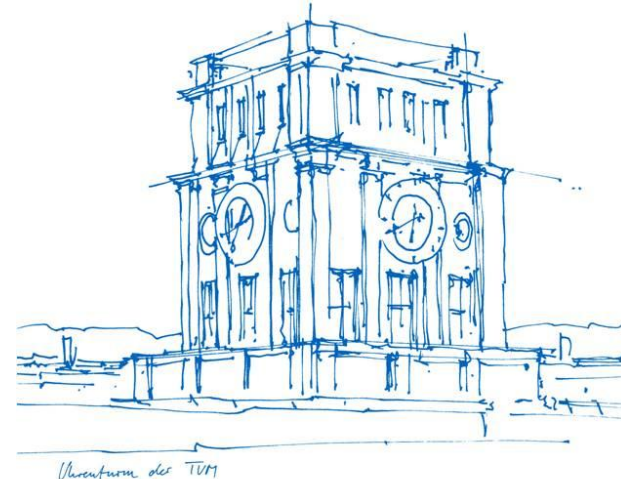
IFAN: An Explainability-Focused Interaction Framework for Humans and NLP Models

Edoardo Mosca, Daryna Dementieva, Tohid E. Ajdari,
Maximilian Kummeth, Kirill Gringauz,
Yutong Zhou, Georg Groh.

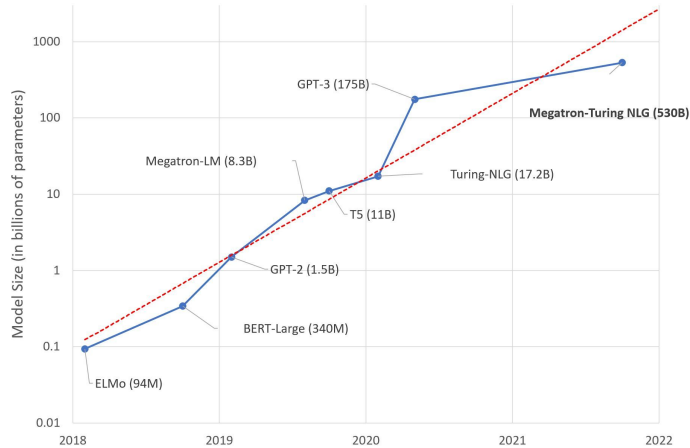
Technical University of Munich, Germany

AAACL 2023

1st-4th October | Bali



The Current Stand



Performance

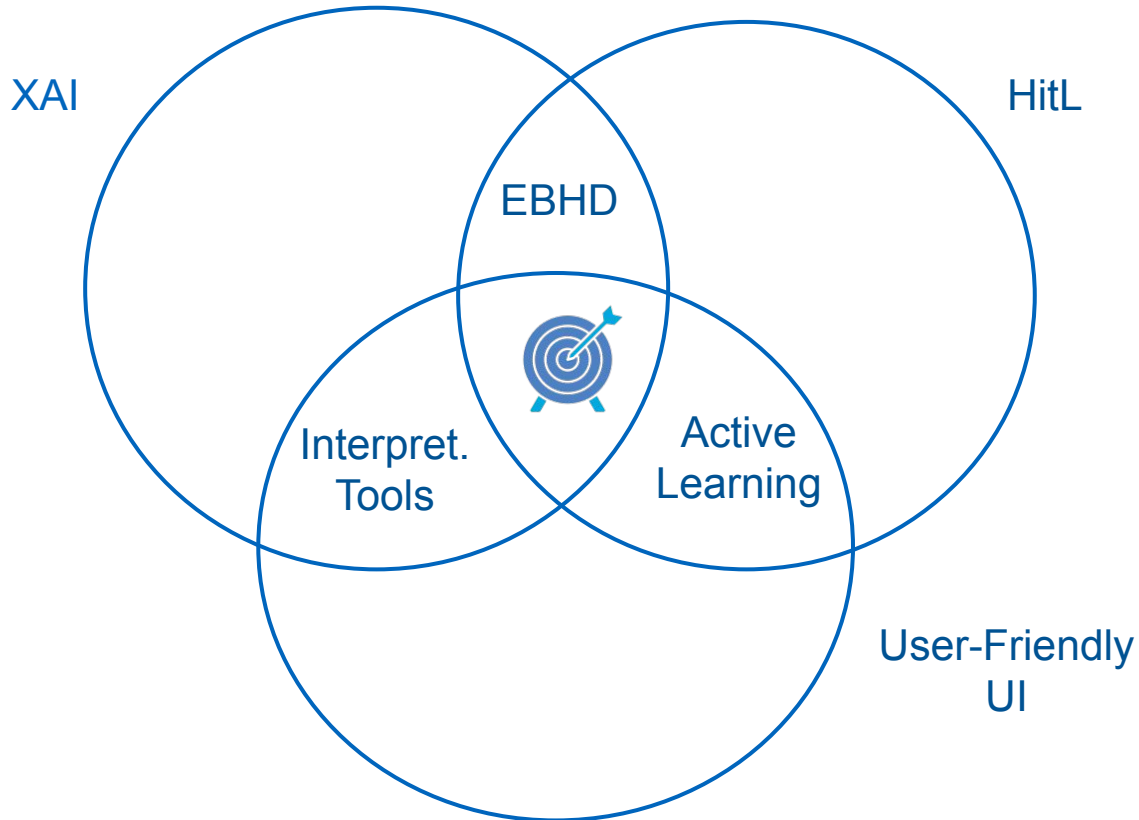


Interpretability

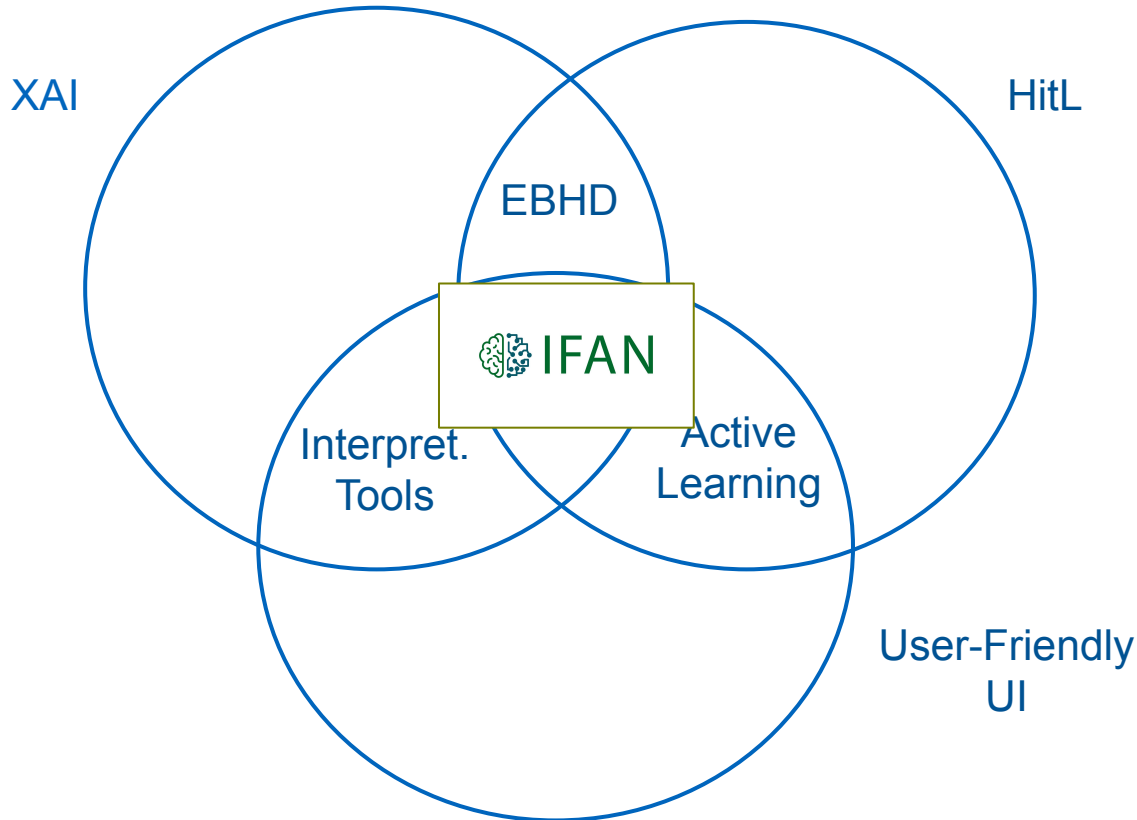


Controllability

The Current Stand



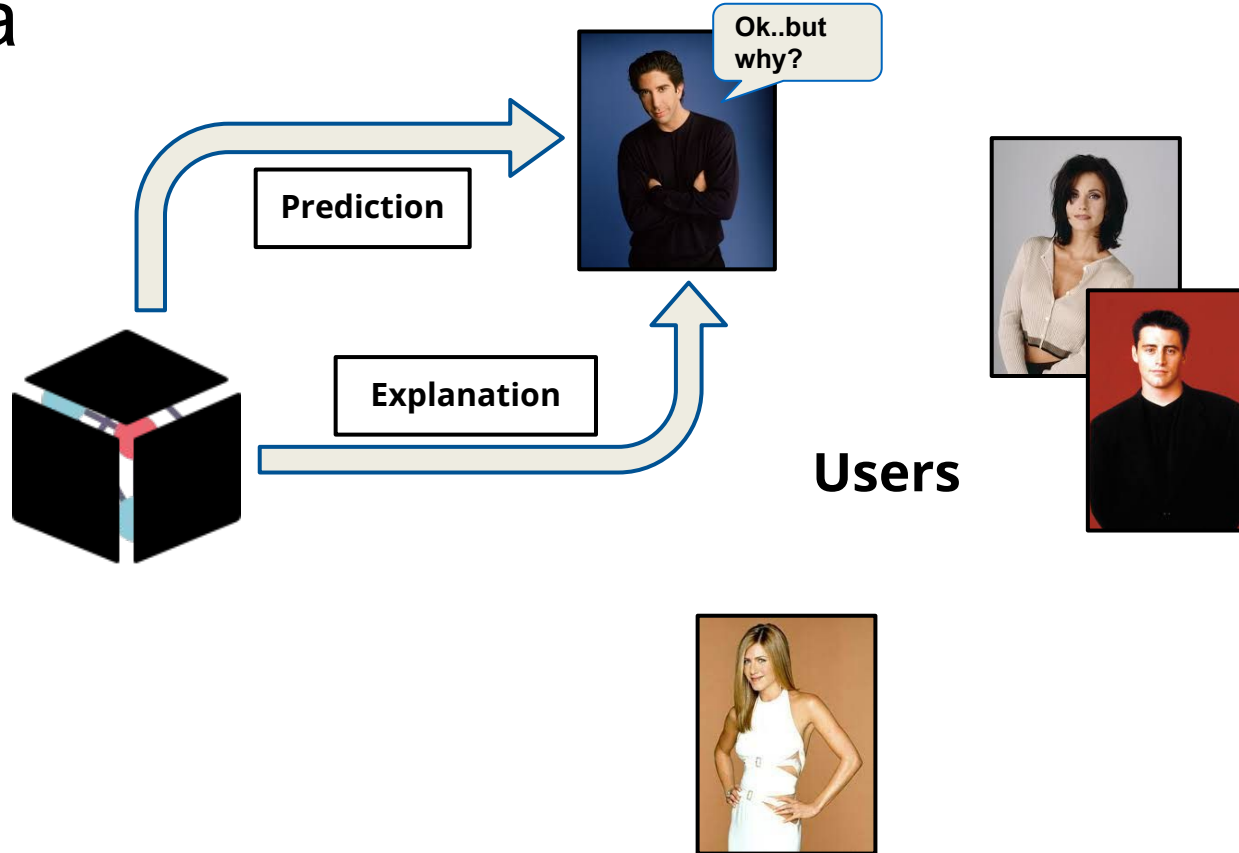
The Current Stand



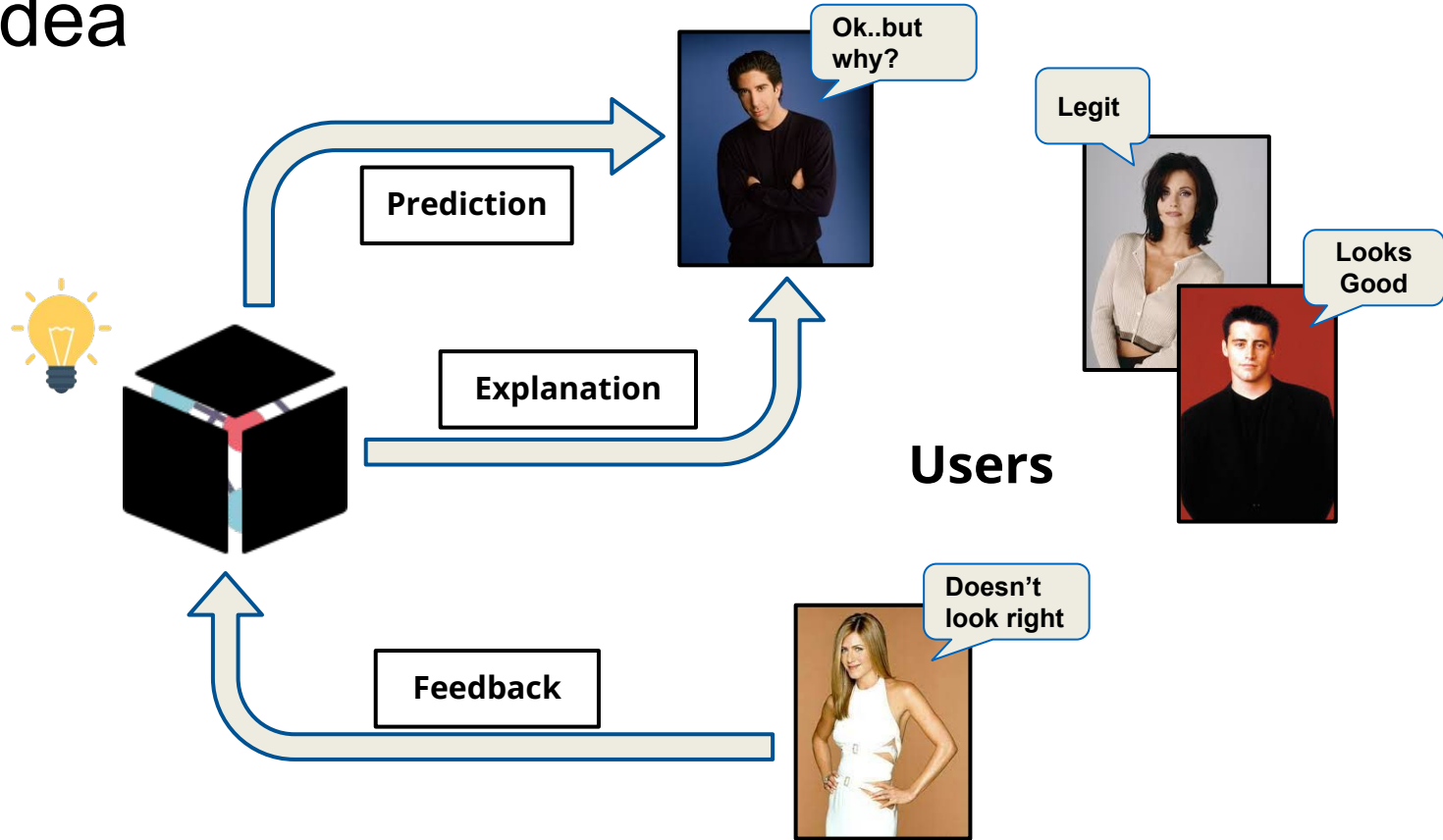
Our contribution

- 1 Propose IFAN: an EBHD Framework for NLP Models**
Users can observe explanations, edit the rationale, give feedback, etc.
- 2 Extend it with a UI and a Management Systems**
Monitor improvement, configure models and user access, etc.
- 3 Test IFAN on a Model Debiasing Task**
We propose feedback rebalancing to contrast model forgetfulness.

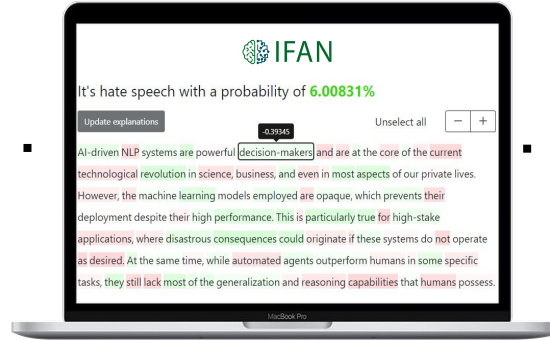
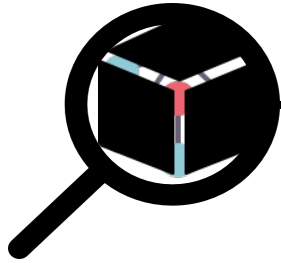
The Idea



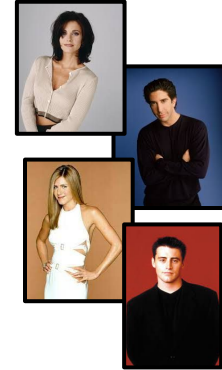
The Idea



The Final Prototype



Users



Manage:
- Models
- Datasets
- Users (Rights)

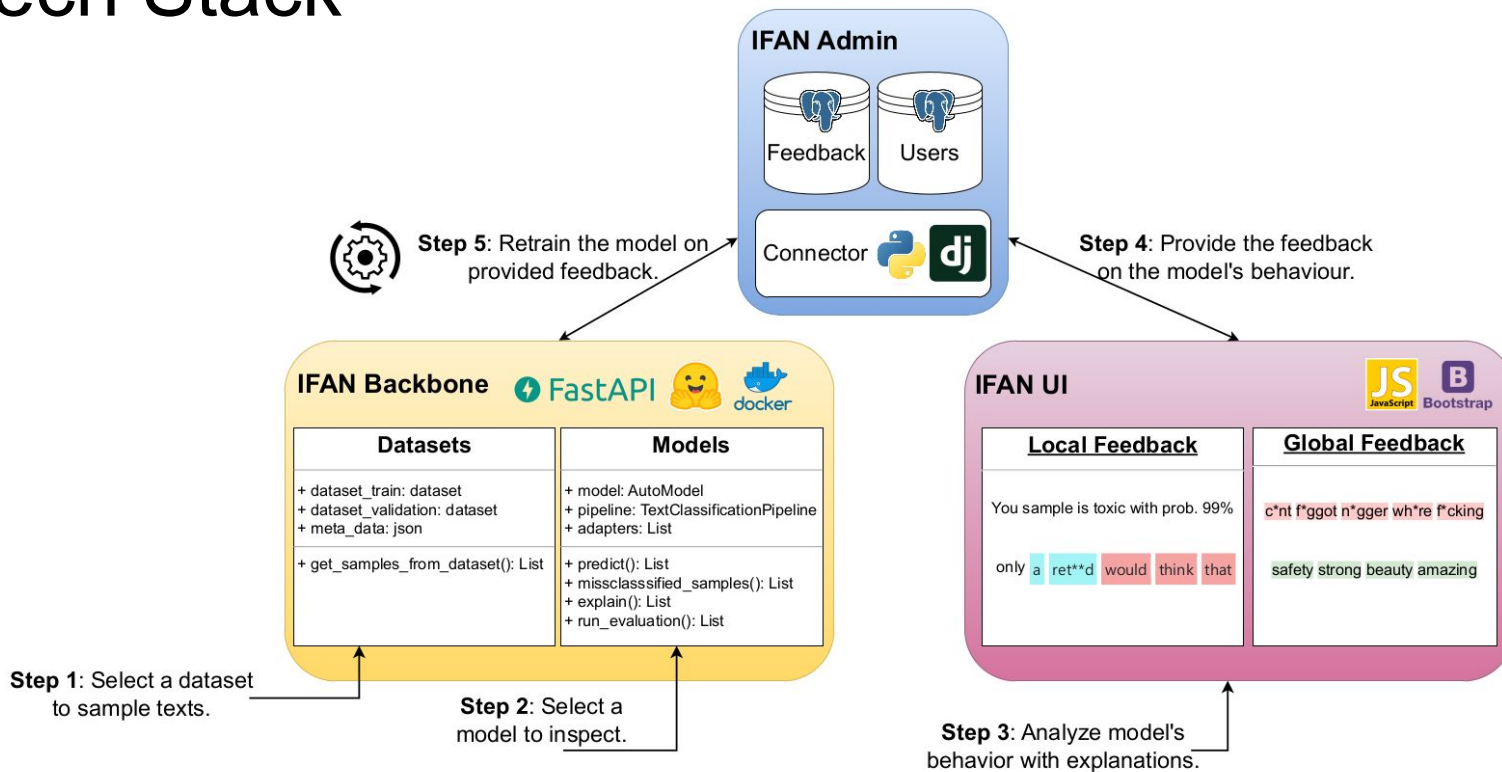
Reports:
- Analyze
- Test
- Monitor

Endpoints:
- Deploy
- Experiment
at Scale



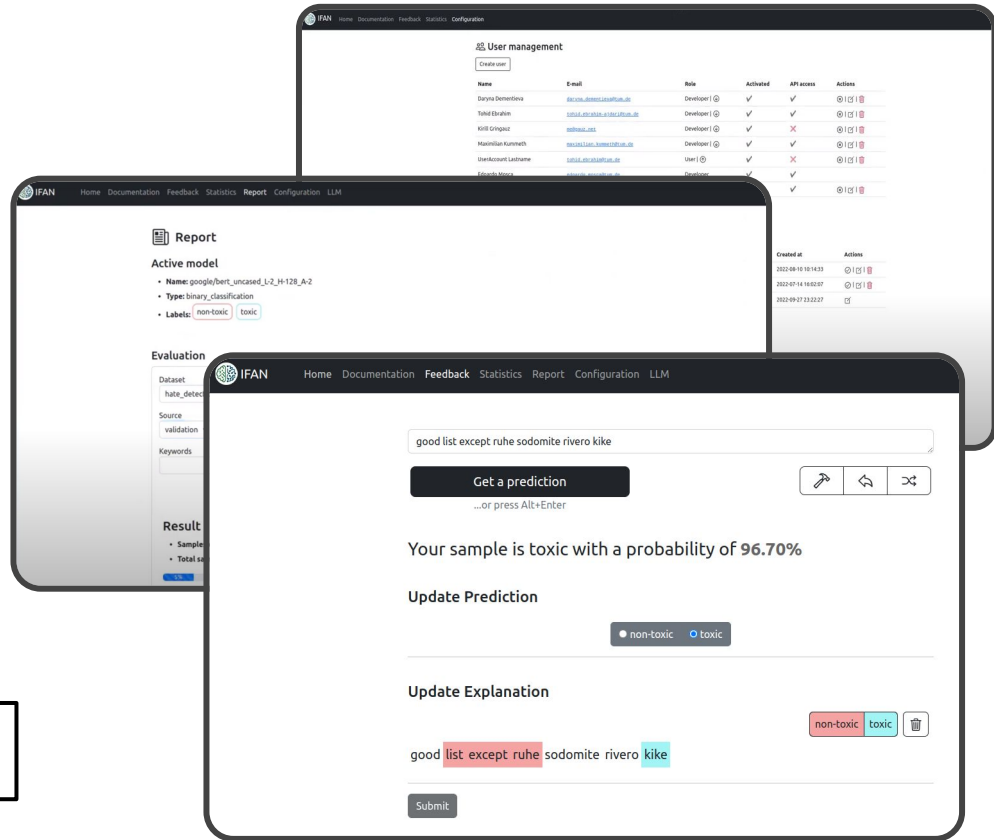
Devs

Tech Stack



- Home
- Documentation
- Feedback
- Report
- Configuration

[To the Demo Video](#)



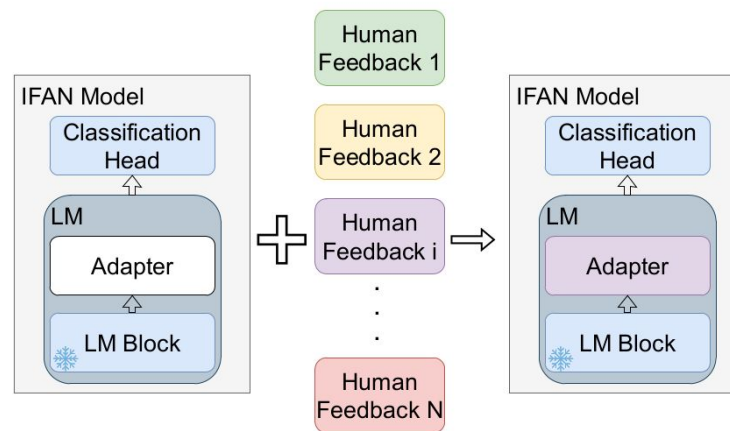
Backbone & API



models		
GET	/models/	Get All Models Api
GET	/models/meta_data	Get Meta Data Of Model Api
PATCH	/models/meta_data	Patch Meta Data Of Model Api
GET	/models/missclassified_samples	Get Misclassified Samples Api
GET	/models/evaluation	Get Evaluation Api
GET	/models/available_adapters	Get Available Adapters Api
POST	/models/train_adapter	Post Update Adapter Api
POST	/models/upload	Upload
datasets		
GET	/datasets/	Get All Datasets Api
GET	/datasets/meta_data	Get Meta Data Of Dataset Api
GET	/datasets/samples	Get Samples Api
GET	/datasets/samples_with_keyword	Get Samples Keyword Api
POST	/datasets/upload	Upload
prediction		
POST	/prediction/	Get Prediction Api
explanation		
POST	/explanation/local/	Post Local Explanation Api
POST	/explanation/local/lime	Post Local Lime Explanation Api
POST	/explanation/local/captum	Post Local Captum Explanation Api

User Roles & Feedback Mechanism

	Dev	Annotator	Unauthorized
Classification & Explanations	✓	✓	✓
Smart Samples Selection	✓	✓	✗
Feedback	✓	✓	✗
Active Configuration	✓	✗	✗
Model Report & Misl. Samples	✓	✗	✗
New Models & Datasets Upload	✓	✗	✗
New Users Creation	✓	✗	✗



Adapters are trained with batches of human-highlighted sample segments.

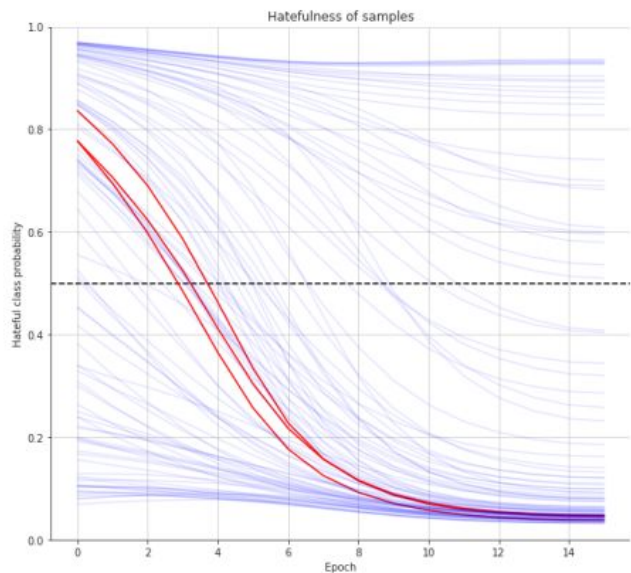
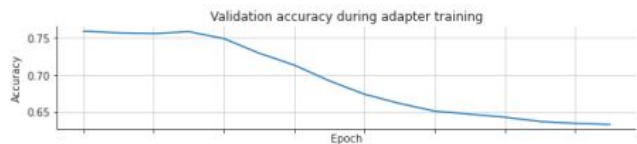
Models love to Forget

- Adapters learn very quickly from feedback (especially on strong signals).
- Models tend to forget what they know, with large drops in performance.
- Adding original samples to feedback batches mitigates the problem (**Rebalancing**).

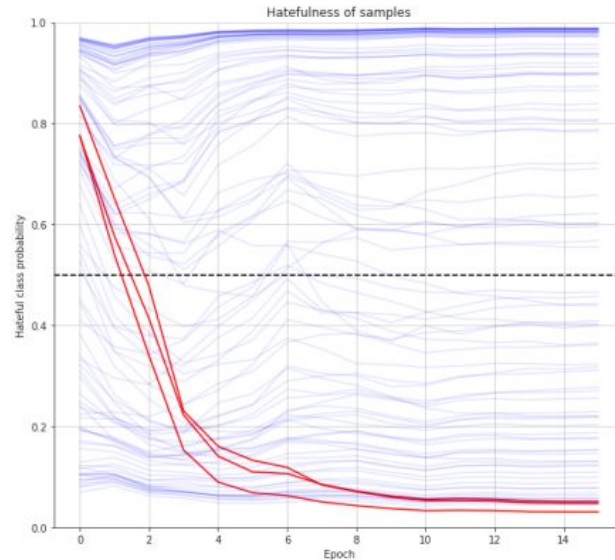
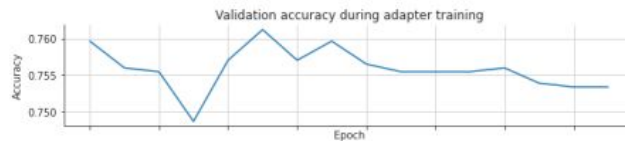
Performance on Use Case

Model	Pr	Re	F1	Pr _J
BERT (baseline)	0.80	0.78	0.79	0.95
<i>Most Confident Missclassified</i>				
BERT+Feedback (non-bal.)	0.34	0.28	0.31	0.82
BERT+Feedback (bal.)	0.78	0.80	0.79	0.97
<i>Least Confident Missclassified</i>				
BERT+Feedback (non-bal.)	0.83	0.73	0.78	0.96
BERT+Feedback (bal.)	0.79	0.78	0.78	0.96

Effects of Rebalancing



Without Rebalancing

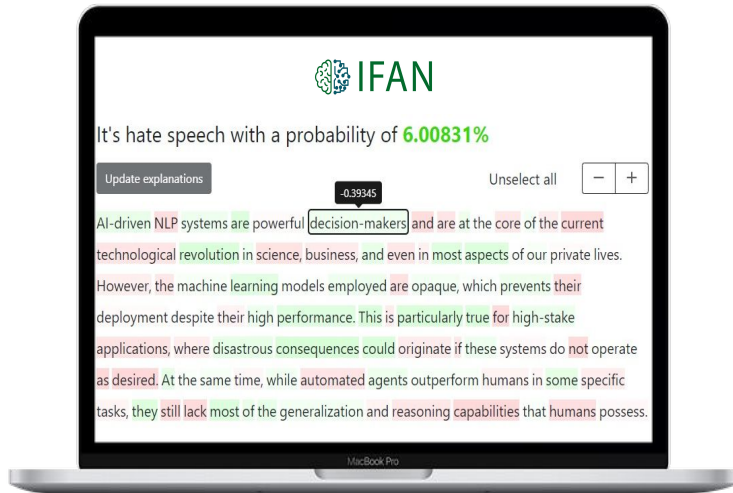


With Rebalancing

Limitations & Takeaways

- Editing explanations can carry human feedback which is beneficial for fixing NLP models.
- As of now we support sentence-to-class and token-to-class tasks (other tasks are work in progress)
- No clear optimal choices of hyperparameters for feedback
- Dealing with a large spectrum of models is hard, but is becoming easier.
- In most cases, your model won't get much better in performance, but is more aligned with humans where it received feedback.
- Rebalancing allows integrating feedback with minimal performance loss.

Thank you!!



Paper



Demo Video

Thank you!!

Daryna Dementieva

Tohid E. Ajdari,

Maximilian Kummeth

Kirill Gringauz

Yutong Zhou

Georg Groh



Edoardo Mosca
edoardo.mosca@tum.de



IFAN