Rethinking Usability Heuristics for Modern Biomedical Interfaces

Stefan Röhrl, M.Sc. DigitalWorld 2023 Congress 16th International Conference on Advances in Computer-Human Interactions

Cellface 🖂 cellface@cit.tum.de 🌐 wiki.tum.de/display/cellface 🗍 Stefan Röhrl, M.Sc. 🏦 Technical University of Munich 🆗 ACHI 2023



Holographic Cell Imaging

Stained Light Microscope Image



https://en.wikipedia.org/wiki/Agranulocyte#/media/File:Monocyte_no_vacuoles.JPG

😂 CelliFace 🖂 cellface@cit.tum.de 🌐 wiki.tum.de/display/cellface 🧏 Stefan Röhrl, M.Sc. 🏦 Technical University of Munich 🍟 ACHI 2023

Holographic Phase Image



User Interface Development Process



CellFace 🖂 cellface@cit.tum.de 🌐 wiki.tum.de/display/cellface 🧏 Stefan Röhrl, M.Sc. 🏛 Technical University of Munich 🖗 ACHI 2023

Speeding up Formative Evaluation



Heuristic Evaluation by Experts

Method Nielsen (1990, 1995)

- Predefined best practice rules
- Team of (five) expert evaluators
- Rate severity of violations [0, 4]
- Overall Summary

Advantages Nielsen (1993), Noyes (1999)

- Time-efficient
- Less Human Resources
- Can be performed by novices and experts equally





Set of (potential) usability problems rated by their expected severity

Goals

Custom Heuristic Rules

- Develop a set of custom heuristics
- Focus on biomedical user group
- Compare their performance against other similar and popular heuristics

User Interface

- Support interdisciplinary research
- Allow interaction of humans and Al
- Fully functioning prototype



Existing Heuristics

Nielsen's heuristics Nielsen (1990,2005) Human-Al interaction heuristics Amershi (2019)

- 1 Simple and natural dialogue
- 2 Speak the user's language
- **3** Minimize the user's memory load
- **4** Consistency and standards
- 5 Feedback
- 6 User control and freedom
- 7 Shortcuts
- 8 Good error messages
- 9 Error prevention
- **10** Help and documentation

- Make clear what the system can do
 Make clear how well the system can do what it can do
 Time services based on context
 Show contextually relevant information
 Match relevant social norms
 Mitigate social biases
 - Support efficient invocation
- 8 Support efficient dismissal
- **9** Support efficient correction

- 10 Scope services when in doubt
 11 Make clear why the system did what it did
 12 Remember recent interactions
 13 Learn from user behavior
 14 Update and adapt cautiously
 15 Encourage granular feedback
- **16** Convey the consequences of user actions
- **17** Provide global controls
- **18** Notify users about changes

CellFace icellface@cit.tum.de wiki.tum.de/display/cellface Stefan Röhrl, M.Sc. Technical University of Munich ACHI 2023

Biomedical AI Heuristics

Requirements

- Generalizing to new interfaces
- Adapted to the domain
- Easy to understand and apply

Approaches

- Adapting existing heuristics
- Reviewing literature of the designated domain
- Expert interviews



Biomedical AI Heuristics



Interface Prototype

/iew 1: Data Setup	View 2: Initialization	View 3: Algorithm Selection
Data Setup Initialization Adjorithm Selection Assisted Training Review Review	Data Setup Data Setup Algorithm Selection Adjoint Selection Adjoint Selection Adjoint Selection Recording (gender 0) Korve labels Recording (gender 2) Active labels Recording (gender 2) Active labels Recording (gender 2) Active labels Recording (gender 4) Active labels	Data Setup Ontalization Agonthm Selection Assisted Training Preview Printh NATVE BAYES RANDOM FOREST K-NEAREST-NEIGHBOURS NEURAL NETWORK Second Segmentation Preview Printh Preview Printh P
EXCE OUSTOMIZE LABELING CLASSES DEXT View 4: Assist © Data Setue @ Initalization @ Agorithm Setue	Image: Construction of the state of the	lection @ Assisted Training @ Review @ Fitch 1.41 you Letts Letters Letters Letters Letters Letters Letters Letters Letters Letters Letters Letters Letters Letters Letters Lett

larald Topfer

NEXT

🛿 🖉 Celll=CCE 🖂 cellface@cit.tum.de 🌐 wiki.tum.de/display/cellface 🕺 Stefan Röhrl, M.Sc. 🏦 Technical University of Munich 🏆 ACHI 2023

Background

Uncertainty: 52.7 %

Comparison of Usability Evaluation Methods

- Evaluate the interface with each Usability Evaluation Method (UEM)
- Conduct a baseline test to create a ground truth of usability problems
- Compare the performance using suitable metrics





Heuristic Evaluation



Human-AI heuristics



Heuristic Evaluation

20

Biomedical-AI Heuristics

Streamline main task Provide full control Orientation Guide attention Provide comparisons Show impact User over system Familiar language Precise language Familiar look Appeal Explain data Explain processing Explain reasoning Strengths Limitations



Heuristics	Predicted Problems
Nielsen's Heuristics	60
Human-Al interaction heuristics	26
Biomedical-Al heuristics	55

- Highest number of violations found by Nielsen's set of heuristics
- Experts stated that Human-Al heuristics were harder to apply

Empirical User Test



- Discovered 75 real usability problems
- Many have already been predicted by heuristics



Lewis (1994), Hartson (2001)

Performance Comparison

Metric	Nielsen's heuristics	Human-Al heuristics	Biomedical-Al heuristics
Precision (Validity)	63.3%	73.1%	85.5%
Recall (Thoroughness)	50.1%	25.3%	62.7%
Average Severity	2.66	2.84	2.74
Weighted Recall	54.0%	28.9%	69.0%
Recall (Severity = 3)	38.9%	27.8%	77.8%
Recall (Severity = 4)	73.3%	46.7%	93.3%

Performance Comparison

Precision in %

Recall in %



Conclusion

Nielsen's heuristics

- Found most genuine problems
- Less efficient AI induced parts
- Reduced precision

Human-AI heuristics

- Worst precision and recall
- Some rules were not effective
- Difficult to apply by experts

Biomedical-AI heuristics

- Satisfying overall performance
- Uncovers most severe problems
- Best adaption to AI in this domain

Overall Remarks

- Demonstrated need for domain specific sets of heuristic rules
- Unweighted recall of 63% is helpful for further improvement
- Further evaluation on different interfaces needed

Thank You





References

S. Amershi, M. Cakmak, W. B. Knox, and T. Kulesza (2014). Power to the people: The role of humans in interactive machine learning. *Al Magazine*, vol. 35, no. 4, pp. 105–120.

S. Amershi, D. Weld, M. Vorvoreanu, A. Fourney, B. Nushi, P. Collisson, J. Suh, S. Iqbal, P. N. Bennett, K. Inkpen et al. (2019). "Guidelines for human-Al interaction. *In Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems,* pp. 1–13.

J. De Fauw, J. R. Ledsam, B. Romera-Paredes, S. Nikolov, N. Tomasev, S. Blackwell, H. Askham, X. Glorot, B. O'Donoghue, D. Visentin et al. (2018). Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature Medicine*, vol. 24, no. 9, pp. 1342–1350.

European Commission (2019). High-Level Expert Group on Artificial Intelligence: Ethics Guidelines for Trustworthy AI.

T. G. Gill (1996). Expert systems usage: Task change and intrinsic motivation. *Management Information Systems Quarterly*, pp. 301–329.

H. R. Hartson, T. S. Andre, and R. C. Williges (2001). Criteria for evaluating usability evaluation methods," International Journal of Human-Computer Interaction, vol. 13, no. 4, pp. 373–410.

References

A. Holzinger, C. Biemann, C. S. Pattichis, and D. B. Kell (2017). What do we need to build explainable ai systems for the medical domain? *arXiv preprint arXiv:1712.09923*.

E. Horvitz (1999). Principles of mixed-initiative user interfaces. *In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp. 159–166.

A. D. Jameson (2009). Understanding and dealing with usability side effects of intelligent processing," *AI Magazine*, vol. 30, no. 4, pp. 23–23.

J.K. Kies, R.C. Williges, and M.B. Rosson (1998). Coordinating computer-supported cooperative work: A review of research issues and strategies. *In Journal of the American Society for Information Science*, vol. 49, no. 9, pp. 776–791.

J. R. Lewis (1994). Sample sizes for usability studies: Additional considerations. *Human factors*, vol. 36, no. 2, pp. 368–378.

H. Lieberman (2009). User interface goals, AI opportunities. AI Magazine, vol. 30, no. 4, pp. 16–22, 2009.

J. Nielsen and R. Molich (1990). Heuristic evaluation of user interfaces. *In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp. 249–256.

References

J. Nielsen and T. K. Landauer (1993). A mathematical model of the finding of usability problems. *In Proceedings of the INTERACT'93 and CHI'93 Conference on Human Factors in Computing Systems*, pp. 206–213.

J. Nielsen (1995). Severity ratings for usability problems. *Papers and Essays*, vol. 54, pp. 1–2.

J. Nielsen (2005). Ten usability heuristics. (accessed 2022.12.17). [Online].

J. Noyes and C. Baber (1999). User-centred design of systems. Springer Science & Business Media.

C. Rzepka and B. Berger (2018). User interaction with Al-enabled systems: a systematic review of is research. *In Thirty Ninth International Conference on Information Systems*, vol. 39, pp. 1–17.

A. Sears (1997). Heuristic walkthroughs: Finding the problems without the noise. *International Journal of Human–Computer Interaction*, vol. 9, no. 3, pp. 213–234.

R. T. Sutton, D. Pincock, D. C. Baumgart, D. C. Sadowski, R. N. Fedorak, and K. I. Kroeker (2020). An overview of clinical decision support systems: benefits, risks, and strategies for success. *Nature Partner Journals Digital Medicine*, vol. 3, no. 1, pp. 1–10.

Biomedical Research AI Heuristics

Streamline main task

Provide full control

Structure

Orientation and Consistency Make the system easy to learn. Focus on the main task that a system was created for and move unimportant side tasks away from the center of attention. The presentation of the main task should be simple to not overwhelm novice users.

Provide global controls and access to important model parameters that may be adjusted to boost the system performance. Complex tasks require fine tuning and advanced users often want to make these changes on their own.

Al systems are often complex which causes users to get lost. Always show the user where he is, what is currently going on and what he can do next. Navigation and interaction items should be easily noticeable and distinguishable.

Biomedical Research Al Heuristics

Guide attention

Interaction

- Provide comparisons
- 6 Show impact
 - User over system

It is easy for the user to lose track what is relevant in a complex AI system. Thus, it is important to guide the user's attention to where his action is needed. Keep him focused on his task and do not throw too many alarms and notifications at him.

Let biomedical users compare among similar data or among parameters when they need to judge an outcome of the system or make a decision on their own.

The user needs to see how his actions and decisions directly influence the system and its performance.

Allow the user to correct errors of the AI efficiently at all times and even turn off the AI completely if needed.

Biomedical Research AI Heuristics



- Familiar language
- Precise language
- 10 Familiar presentation styles
- 11 Appeal

8

9

Use non-technical language if possible. Pay attention to use the correct terminology for medical concepts.

Avoid ambiguous wording for labels and commands that could trigger confusion in the user regarding how to do a task or how something works.

Do not use ways of presentation for the interface that biomedical users do not know from other interfaces they work with and that they cannot interpret.

Al interfaces should have an appealing and premium look. This gives the user a feeling of using a state-of-the-art, high quality and thus trustworthy product.

Biomedical Research AI Heuristics



12 Explain data

- **13** Explain process steps
- **14** Explain reasoning
- **15** Highlight strengths and limitations

The data that a task is based on needs sufficient explanation. The user has to develop an understanding of how it can be interpreted and how it is different from other data he may be familiar with.

There needs to be a high-level explanation for the overall procedure that is performed by the system.

There has to be an explanation why and how the system derived a certain result or prediction.

Show what the strengths and weaknesses of the system are and what expectations are realistic. Highlight situations when the system may be underperforming or is in doubt.

Asymptotical User Testing

- Discovery of usability problems can be modeled as a Poisson process
- Number of detected usability problems reaches an asymptotic maximum
- Test with 21 users with a medical, biological and chemical background (7 students and 14 researchers)
- Worst Case Assumption



Interface Prototype



Human Assisted Labeling Tool

- Biomedical cell segmentation and classification
- Medical user and machine learning algorithm work hand in hand

View 1: Data Setup



🔋 CelliFace 🖂 cellface@cit.tum.de 🌐 wiki.tum.de/display/cellface 🖗 Stefan Röhrl, M.Sc. 🏦 Technical University of Munich 🍟 ACHI 2023

View 2: Initialization



CellFace 🖂 cellface@cit.tum.de 🌐 wiki.tum.de/display/cellface 🖟 Stefan Röhrl, M.Sc. 🏛 Technical University of Munich 🖗 ACHI 2023

View 3: Algorithm Selection

	NAIVE BAYES RAND	OM FOREST K-NEAREST-NEIG	GHBOURS NEURAL NETWORK >	
1. Three	shold Segmentation	2. Feature Extraction	3. Random Forest Classification	
Binariza 50	tion Threshold on on	Morphological Features Area, Perimeter, Circularity,	Number of trees 100	
Minimu 30	m Cell Size px on	Optical Height	Maximum tree depth	
AD	VANCED SETTINGS	ADVANCED SETTINGS	ADVANCED SETTINGS	
Randor	nly builds an ensemble of dec	ision trees. Each decision tree consist	s of layers in which the data is split into groups	

CellFace 🖂 cellface@cit.tum.de 🌐 wiki.tum.de/display/cellface 🖟 Stefan Röhrl, M.Sc. 🏦 Technical University of Munich 🖗 ACHI 2023

View 4: Assisted Training



CellFace 🖂 cellface@cit.tum.de 🌐 wiki.tum.de/display/cellface 🖟 Stefan Röhrl, M.Sc. 🏛 Technical University of Munich 🖗 ACHI 2023

View 5: Review

	0 9 18 27 36 45 54 63 72 81	1.41 seconds/label	49 LABELS	
	Labeling Progress	Labeling Speed	All Labels	
	Number of assigned labels over time	Average label time	Click to display all labels	
	Your sample	Name Speet	d (s/l) Date	
		You 1.41	Today	
	Leucocyte Aggregate Other Platelet	Barry Allen 1.97 Harald Töpfer 3.32	2020-05-18 2020-05-14	
	Label Components	Leader	board	
BACK	Compare the composition to similar datasets	Ton labeling speeds ov	ver the last 6 months	NEXT