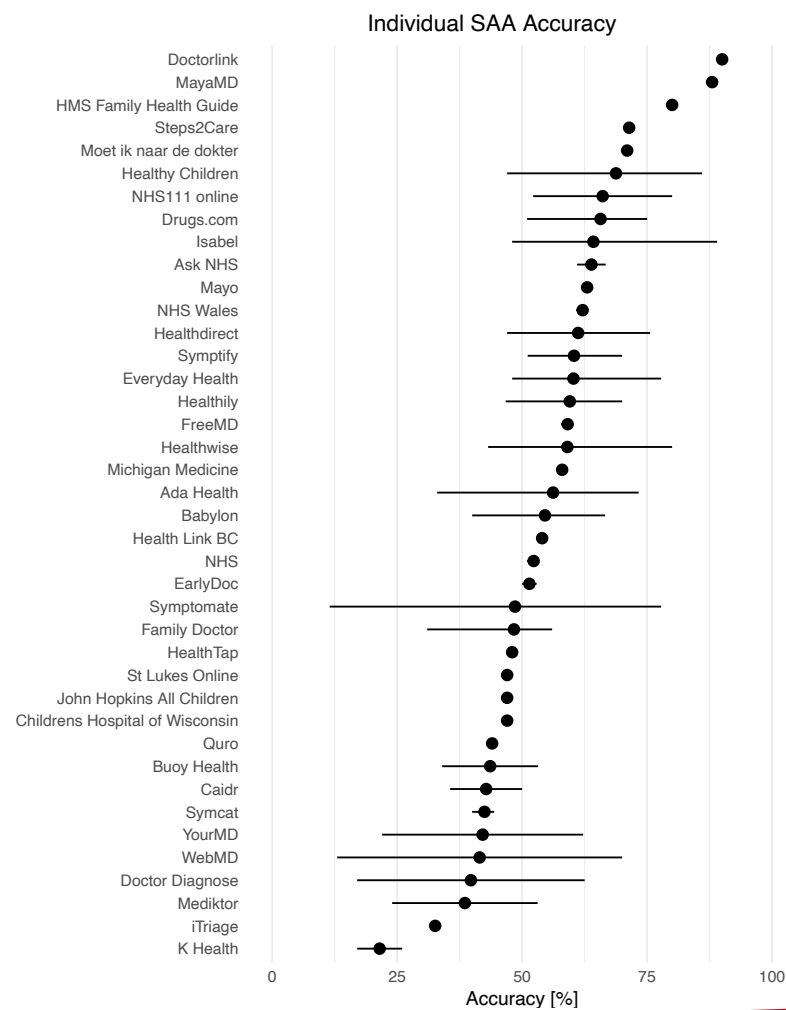


# Vergleich von Ersteinschätzungssystemen

Dr. Dr. Marvin Kopka | Urgent Care Conference | 24.06.2026

Fachgebiet Arbeitswissenschaft, Technische Universität Berlin

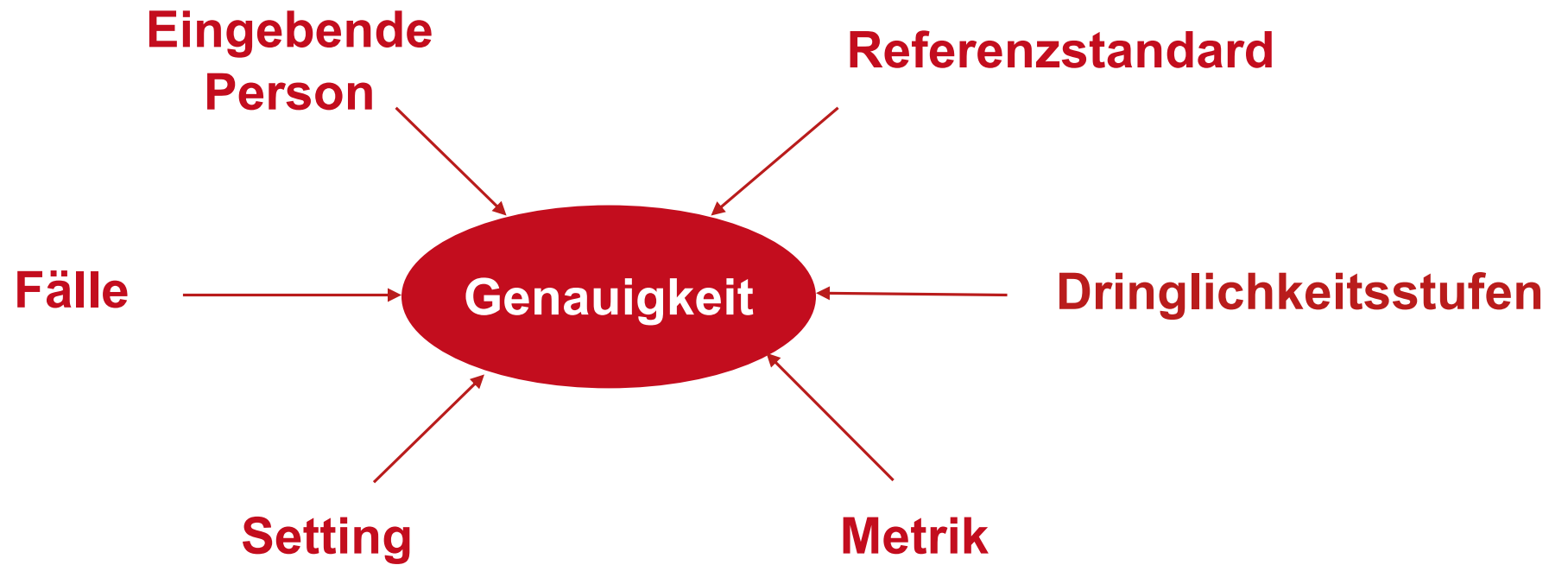
# Genauigkeit von Systemen



# Genauigkeit von Systemen



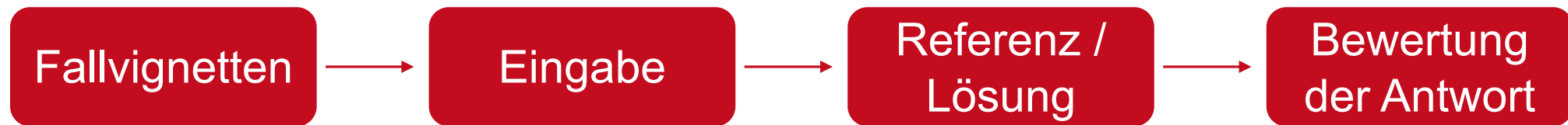
# Genauigkeit?



# Ebenen der Evaluation



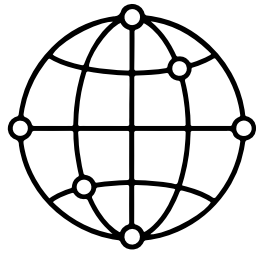
# Vorgehen bei Fallvignettenstudien



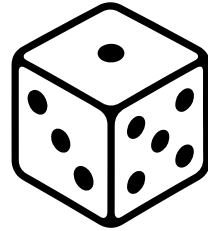
## Probleme mit Fallvignetten

- Beschreibung von echten Fällen?
- Beschreibung von Ärzten oder von Patienten selbst?
- Menge an Informationen? Laborergebnisse?
- Prävalenz von Symptomen repräsentativ?
- Anzahl an Vignetten?

# Lösungen für **Fallvignetten** - Repräsentative Patientenfälle



r/AskDocs  
(n = 8.794)



Zufallsziehung  
& Bewertung



Vorläufiges  
Vignettenset  
(n = 45)

statistische  
Kriterien

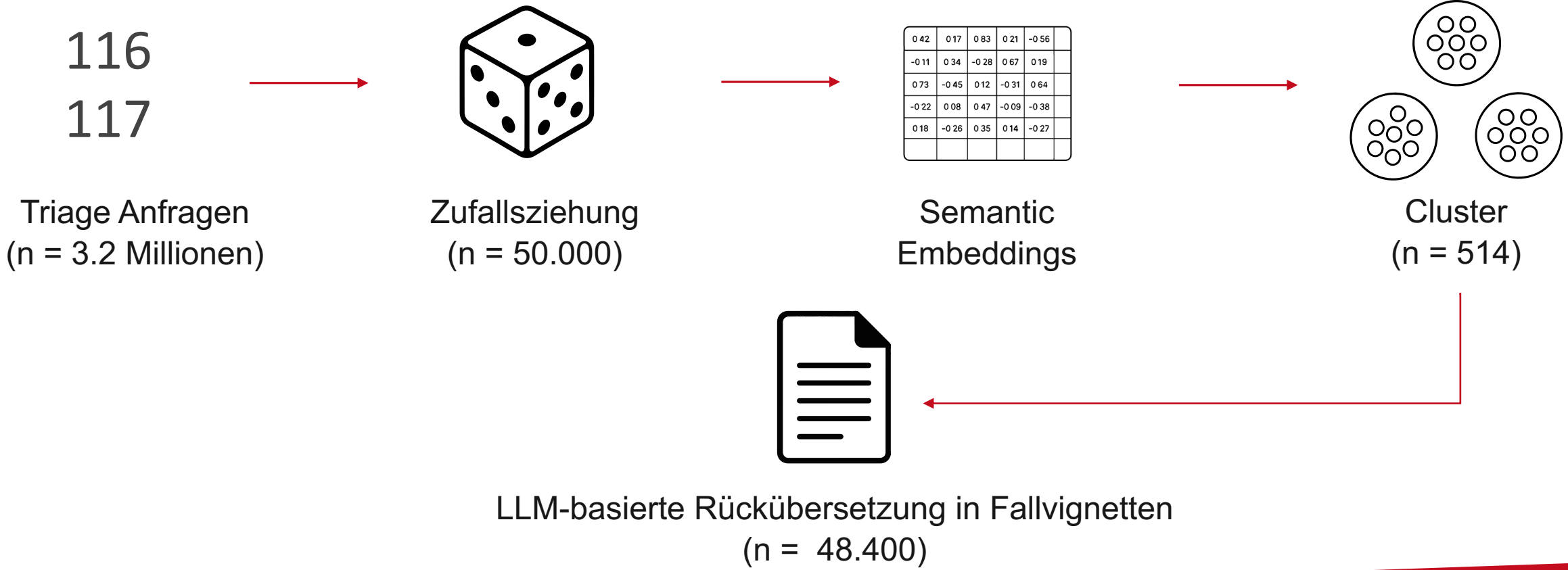


Finales  
Vignettenset  
(n = 27)



Prävalenz

# Lösungen für **Fallvignetten** - Repräsentative Nutzerfälle



Probleme mit

Eingabe

- Wer gibt Fälle ein? Laien? Experten?
- Welche Annahmen werden getroffen?
- Fehlende Information durch klinische Expertise gefüllt?
- Wie viele Personen sollten Fälle eingeben?



16:45 H  
39 Grad, sollte ich damit zum  
Themen gehen?

I'll check reliable medical guidance and give you a clear recommendation, including warning signs and when to seek urgent care.

Thought for 2m 29s >

Ja – mit 39 °C Fieber seit 4 Tagen plus Halsschmerzen solltest du dich heute noch ärztlich abklären lassen. Das ist nicht mehr etwas, das ich einfach weiter beobachten würde. Länger anhaltendes Fieber oder Temperaturen über 39 °C können ein Hinweis sein, dass sich eine Infektion verschlimmert oder behandlungsbedürftig ist.

Da heute Samstag ist: Wenn deine Hausarztpraxis zu ist, ruf den ärztlichen Bereitschaftsdienst 116117 an oder nutze eine Bereitschaftspraxis. Die 116117 kann Beschwerden gedanklich aufnehmen und dich normalerweise in eine Praxis empfehlen, aber nicht bis zum nächsten Morgen sollten.

Sofort 112 statt 116117, wenn eines der folgenden Anzeichen vorliegt:

+ Ask ChatGPT

# Lösungen für **Eingabe** - Adjustierte Metriken & Eingabeprotokoll



- Eingabe von mind. 3 Personen
- Meta-Genauigkeit: Genauigkeit bei der Eingabe von Symptomen
- Standardisierte Eingabeprotokolle



Chief complaint(s)	<ul style="list-style-type: none"> <li>- Copy and paste the prescribed chief complaint(s) word by word (if only one symptom can be added at the time, add one by one)</li> <li>- If none of the symptoms are accepted by OSC, stop the vignette entry</li> <li>- In case of multiple chief complaints, if at least one of the symptoms is successfully inputted, continue the consultation.</li> <li>- If not all the symptoms are successfully inputted but there is an option to add more, attempt to do so. If unsuccessful, continue with the consultation without adding them.</li> </ul>
Demographics	<ul style="list-style-type: none"> <li>- Copy over/input age and gender of the persona of the vignette</li> </ul>
Instruction for answering questions during the question-answer part of the consultation	<ul style="list-style-type: none"> <li>- Only select symptoms that are mentioned in the vignette</li> <li>- Do not to select any symptom that is not described in the vignette even if you feel it might be missing from the vignette</li> <li>- Caveats:                             <ul style="list-style-type: none"> <li>- If synonyms offered for symptom described in the vignette, select it e.g. "abdominal pain" for "tummy pain"</li> <li>- If a wider logical category of a symptom is offered, select it e.g. if "knee pain" is in the vignette, then select "leg pain" if offered</li> <li>- If the sentence in the vignette is worded differently but has the same meaning, select it.</li> <li>- If a sentence in a vignette indirectly could lead to a symptom or vice versa, then do not select that symptom. E.g. if someone is 'waking up at night with pain', but the symptom checker asks if the 'pain is severe' - do not select that symptom as it is unknown whether the symptom is severe, just they are waking up at night with it. However, if the vignette states that someone is 'crying with the pain' then you can select it, as it may be understood as a <u>synonym</u> of 'pain is severe'</li> <li>- If further characteristics of a confirmed symptom is asked in one question with multiple answers, and there is an option to decline all of them, do so. If this is not possible but there is an option not to answer e.g. "I do not know"/"Not prefer to say", choose that, otherwise, stop the consultation and mark the case as "INCOMPLETE"</li> </ul> </li> <li>- If some of the chief complaint(s) are not recognised but asked later during the consultation, confirm them</li> </ul>
Documentation	<ul style="list-style-type: none"> <li>- Free text entry and other entered symptoms via drop-down</li> <li>- Confirmed chief complaint(s)</li> <li>- Outcome conditions</li> <li>- Triage advice</li> <li>- All the data points (including confirmed and declined symptoms) - depending on symptom checker, if possible copy the consultation report with as much detail as possible</li> </ul>

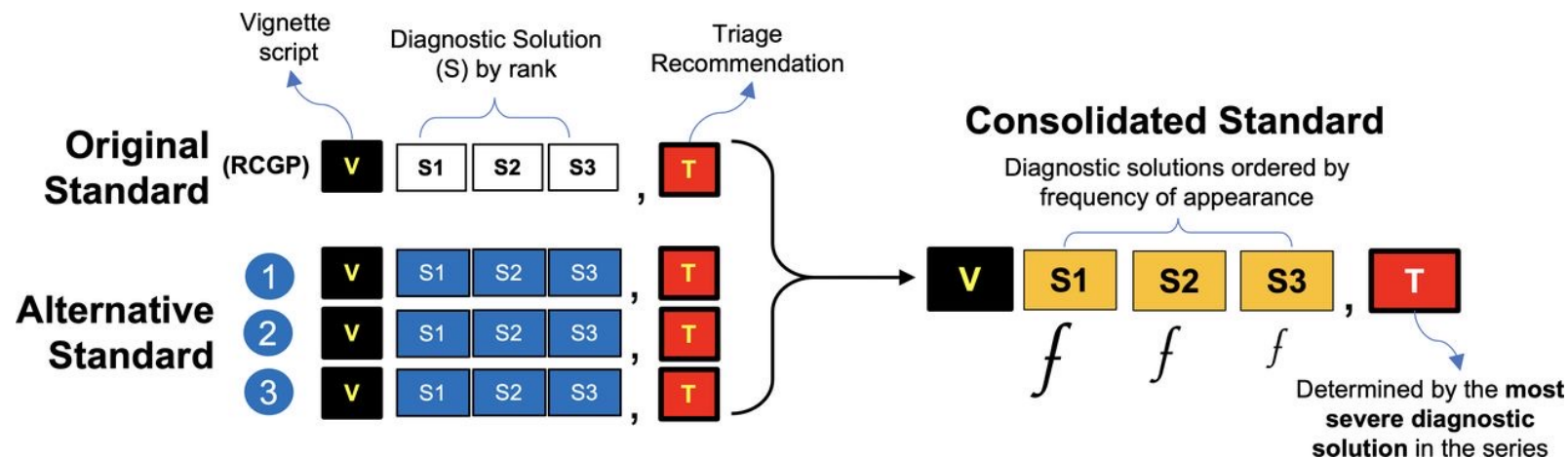
## Lösungen für **Eingabe** - LLMs als Patientensimulation

- „Patient Agent“
- LLM erhält Vignette
- Antwortet nur auf die gestellte Frage auf Basis von Vignettendaten
- Vermeidet klinische Ausdrucksweise
- Füllt fehlende Informationen mit plausiblen Daten

- Welches Verfahren?
- Finale Diagnose und Empfehlung? Bewertung durch Arzt? Ärzte? Ärztepanel?

# Lösungen für **Referenz / Lösung** - Ärztepanel

- Mehrere Ärztepanel bewerten Fälle unabhängig voneinander
- Lösungen werden von weiterem Panel aggregiert



# Lösungen für **Referenz / Lösung** - Richtige & plausible Lösung

- Festlegung einer **richtigen** Lösung
- Festlegung weiterer **plausibler** Lösungen
- Genauigkeit für richtige und plausible Lösungen

- Reicht Genauigkeit?
- Welche weiteren Metriken sollten berichtet werden?
- Wie wird Empfehlung von System den richtigen Lösungsstufen zugeordnet?

# Lösungen für **Bewertung der Antwort** - §360b GeDIG

- Genauigkeit
- Sensitivität (-> Positive Predictive Value?)
- Spezifität (-> Negative Predictive Value?)
- Weitere Metriken (Laienverständlichkeit, Informationssicherheit, ...)

# Lösungen für **Bewertung der Antwort** - symptomcheckR



- Systematische Übersicht zu berichteten Metriken
- Genauigkeit, Genauigkeit pro Level, Sicherheit, Vollständigkeit, Übertriage, Untertriage, Capability Comparison Score
- R-Package zur einfachen Bestimmung der Metriken

## Bewertung der Antwort

# - symptomcheckR



## Checklist – Symptom Checker Accuracy Reporting (SCARF)



### How to Evaluate the Accuracy of Symptom Checkers and Diagnostic Decision Support Systems: Symptom Checker Accuracy Reporting Framework (SCARF)

[Marvin Kopka<sup>1</sup>](#) ; [Markus A Feufel<sup>1</sup>](#) 

- 19 Themen: Vignetten, Referenzstandard, Eingabe, Metriken, Analyse, ...
- Standardisierung von Benchmarkingstudien
- Transparenz

# Checklist – Symptom Checker Accuracy Reporting (SCARF)

Topic	Item Number	Item Description	Page Number
<b>Title &amp; Abstract</b>			
Title	1	Title should indicate that the study evaluates a symptom checker or diagnostic decision support system	
Abstract	2	Summary of evaluation objective, methods, results, and conclusions	
<b>Introduction</b>			
Background and Objectives	3a	State the intended use case of the symptom checker (e.g., self-triage, emergency care triage)	
	3b	Define the target population to which findings are intended to generalize	
<b>Methods</b>			
Case Vignettes	4a	Describe the source of vignettes (e.g., medical education textbooks, patient records, patients' descriptions, case studies, fictitious)	
	4b	Report the sampling frame and rationale (e.g., which conditions, prevalence data, population-level statistics informed vignette selection or creation)	
	4c	Report how statistical representativeness of vignettes was ensured (i.e., representative of prevalence within the sampling frame)	
	4d	Report how content representativeness of vignettes was ensured (e.g., real cases, cases derived from patient records)	
	4e	Report whether atypical cases were included or excluded	
	4f	Specify the number of vignettes and provide a rationale (e.g., power analysis, feasibility)	
	4g	Describe refinement and selection procedures (e.g., test-theoretical metrics such as item-total correlations or item difficulty indices)	
	4h	State whether vignette content was lay-friendly or phrased for clinicians	
Gold Standard Assignment	5a	Describe how the reference standard was established	
	5b	Report the number and background of experts involved	
	5c	Explain how symptom checker outputs were mapped to triage categories	

## Offene Fragen

- Wie viele Fälle sind sinnvoll und machbar?
- Qualität von LLM Patientensimulatoren?
- Genauigkeit von Referenzstandard?
- **Sinnvolle** Metriken?

# Kontakt



[Linkedin.com/in/marvin-kopka](https://www.linkedin.com/in/marvin-kopka)



[@kopka\\_sci](https://twitter.com/kopka_sci)



[mail@marvinkopka.com](mailto:mail@marvinkopka.com)