

# An evaluation framework for unconstrained conversational agents in healthcare: a scoping review

Hang Ding<sup>1</sup>, Joshua Simmich<sup>1</sup>, Atiyeh Vaezipour<sup>1</sup>, Nicole Andrews<sup>1,2,3</sup>, Trevor Russell<sup>1</sup>

<sup>1</sup> RECOVER Injury Research Centre, Surgical Treatment and Rehabilitation Services, Herston, The University of Queensland, QLD, Australia  
<sup>2</sup> Tess Cramond Pain and Research Centre, Metro North Hospital and Health Service, Herston, QLD, Australia  
<sup>3</sup> Occupational Therapy Department, The Royal Brisbane and Women's Hospital, Herston. QLD, Australia



### Introduction

Conversational agents (CAs), also known as chatbots or virtual assistants, are software programs designed to imitate human conversations to engage with users.

Over the past decade, CAs have been increasingly studied to address care burdens and service needs in healthcare delivery. For example, studies have demonstrated the potential to use CA-enabled care programs to assist in triaging people with unurgent health conditions, supporting inpatient care, providing post-discharge follow-ups, and self-managing mental health and chronic diseases.

To use CAs safely and effectively, rigorous evaluations are essential but challenging to achieve. Although reviews have provided some technical metrics and outcome measures, selection of these variables for individual studies remains challenging. To achieve robust and effective outcomes, evaluation frameworks are often needed to help understand essential design, outcome measures, and evaluation targets at different trial stages. However, such frameworks remain absent.

**Our aim: To synthesize existing knowledge and outline a framework for evaluating CAs in healthcare.**

### Method

We conducted a scoping review according to the PRISMA Extension for Scoping Reviews.

We searched CINAHL, Medline, Scopus, Embase

and IEEE Xplore, focusing on CAs which were unconstrained by predefined answers or options. We reviewed study designs, categorised outcome measures, and finally outlined an evaluation framework which can be used to evaluate CAs in healthcare.

To synthesise the framework, we extracted study designs and outcome measures, nested them within well-recognised categories, and mapped the results on an established overarching framework for digital health evaluations.

### Results

The search identified 1553 articles, of which 43 studies were included in the review.

We identified 23 quasi-experimental studies, nine randomised controlled trials, four observational studies, and seven research test-based studies

A total of 175 outcome measures which were used in the reviewed studies were nested into seven categories: 1) functionality, 2) clinical / health outcomes, 3) user experience, 4) costs and cost benefits, 5) safety and information quality, 6) usage, adherence and uptake, and 7) user characteristics for implementation research. We consolidated the framework in the figure below, which shows the outcome measures across four evaluation stages: I) feasibility/usability, II) efficacy, III) effectiveness, and IV) implementation science.

### Discussion

In this review, we found that many studies evaluated CAs using predefined questions, unable to reflect overall CA performances. New strategies to comprehensively evaluate CAs are needed.

We highlighted that some evaluations involving researchers and clinicians are essential, especially for evaluating information quality and risks of interventions, because normal users or participants often have insufficient knowledge to judge the quality or safety of CA-based interventions.

The outcome measures of functionality, user experience, safety and information quality are often diverse and inconsistent. Validated questionnaires and algorithms are needed to achieve robust and effective evaluations.

We identified several measures to be used in future research, such as stability, consistency, standard compliance, and measures to reduce health inequities.

This study is limited to using a single evaluation framework to map the investigation results.

### Conclusion

This systematic review presents a consolidated evaluation framework which can be used to evaluate the performance of CAs in healthcare.

**Figure: The consolidated framework for evaluating conversational agents in healthcare. The framework demonstrates the studies (n=43) and their outcome measures at four major evaluation stages of an established practical guide, named Monitoring and Evaluating Digital Health Intervention, the World Health Organization. Note: Essential measures at different stages, which we identified, are mark by a light blue. \* denotes the measures which we proposed to be included in future studies.**

	Stage →	1. Feasibility & Usability →		2. Efficacy →		3. Effectiveness →		4. Implementation Science	
WHO Digital Health	Brief description	Feasibility: The ability to work as intended. Usability: The degree of a system being used to achieve specified goals in a specified context of use.		Efficacy: The ability to achieve the intended results in a research setting or trial.		Effectiveness: The ability to achieve the intended results in a real application (non-research setting).		Implementation science: To assess the uptake, integration and sustainability of evidence-based digital health interventions for a given context, including policies and practices.	
	Evaluation targets	♦ Stability (system uptime/failure rates) ♦ Performance consistency ♦ Standards adherence (terminology, interoperability, security)		♦ User satisfaction ♦ Workflow “fit” ♦ Learning curve (design) ♦ Cognitive performance / errors ♦ Reliability		♦ Changes in care processes (time) ♦ Changes in outcomes (system performance / health)		♦ Changes in process, outcome, coverage, and costs ♦ Total cost of implementation, and health impact ♦ Error rates ♦ Learning curve of users ♦ Changes in policy, practices attributable to system ♦ Adaptability and extendibility to new use-cases	
Studies Reviewed and Outcome Measures at Four Major Evaluation Stages Aligned with the WHO Guide	Studies included in the review	Almusharraf 2020 <sup>1</sup> Gabrielli 2020 <sup>2</sup> Gaffney 2020 <sup>3</sup> Boczar 2020 <sup>4</sup> Bonnevie 2020 <sup>5</sup> Denecke 2020 <sup>6</sup> Rehman 2020 <sup>7</sup> Stephens 2019 <sup>8</sup>	Bibault 2019 <sup>9</sup> Park 2019 <sup>10</sup> Suganuma 2018 <sup>11</sup> Bickmore 2018 <sup>12</sup> Philip 2014 <sup>13</sup> Yasavur 2014 <sup>14</sup> Rhee 2014 <sup>15</sup>	Davis 2020 <sup>16</sup> Polignano 2020 <sup>17</sup> Caballer 2020 <sup>18</sup> Maher 2020 <sup>19</sup> Lee 2020 <sup>20</sup> Bennion 2020 <sup>21</sup> Bian 2020 <sup>22</sup>	Auriacombe 2018 <sup>23</sup> Fulmer 2018 <sup>24</sup> Fitzpatrick 2017 <sup>25</sup> Friederichs 2014 <sup>26</sup> Harless 2009 <sup>27</sup>	Maeda 2020 <sup>28</sup> Dosovitsky 2020 <sup>29</sup> Chaix 2019 <sup>30</sup> Bott 2019 <sup>31</sup> Perski 2019 <sup>32</sup>	Philip 2017 <sup>33</sup> Ly 2017 <sup>34</sup> Crutzen 2010 <sup>35</sup>	Yang 2021 <sup>36</sup> Fan 2021 <sup>37</sup> Schindler-Ruwisch 2020 <sup>38</sup> Kocaballi 2020 <sup>39</sup> Ferrand 2020 <sup>40</sup>	Nobles 2020 <sup>41</sup> Boyd 2018 <sup>42</sup> Miner 2016 <sup>43</sup>
	Study designs	Single-arm studies <sup>1-8,10,12-15</sup> Randomized controlled trials <sup>9</sup> Two-arm quasi-experimental study <sup>11</sup>		Single-arm studies <sup>16-20,23,27</sup> Case-control study <sup>22</sup> Randomized controlled trials <sup>21,24-26</sup>		Cross-over study <sup>33</sup> Case-control study <sup>31</sup> Cross-sectional study <sup>29,30,35</sup> Randomized controlled trials <sup>28,32,34</sup>		Research test <sup>36,38-43</sup> Cross-sectional study <sup>37</sup>	
	User characteristics	- User characteristics <sup>5</sup>		-		- User characteristics <sup>30,35</sup>		- User characteristics <sup>37</sup> - Users in geographic regions * - Gender, equity and rights – to reduce health inequities *	
	Usage, adherence and uptake	- Usage <sup>3,5,8</sup> - Uptake <sup>5,15</sup>		- Usage <sup>16,19,21,22,24</sup> - Uptake <sup>26</sup> - Adherence <sup>16,22,26</sup>		- Usage <sup>29,30,32,35</sup> - Uptake * - Adherence <sup>34</sup>		- Usage <sup>37</sup> - Uptake * - Adherence <sup>37</sup>	
	Costs and cost benefits	- Costs <sup>5</sup>		- Cost effectiveness <sup>22</sup>		- Costs * - Cost effectiveness *		- Implementation costs *	
	Clinical / health outcomes	- Knowledge and skills <sup>3</sup> - Health wellbeing and issues <sup>11</sup> - Psychological / mental health <sup>3,11</sup> - Clinical assessment performance <sup>13</sup> - Behavioral modification and risk factors <sup>8</sup>		- Knowledge and skills <sup>21,27</sup> - Psychological / mental health <sup>21,24,25</sup> - Clinical assessment performance <sup>18,20,23</sup> - Behavioral modification and risk factors <sup>16,19,26</sup>		- Knowledge and skills <sup>28</sup> - Health wellbeing and issues <sup>31,34</sup> - Psychological / mental health <sup>28,31,34</sup> - Clinical assessment performance <sup>33</sup> - Behavioral modification and risk factors <sup>28,30,32</sup>		- The effectiveness of the approved intervention in less controlled environment *	
	User experience	- Usability <sup>15</sup> - Feasibility <sup>2,8</sup> - Usefulness <sup>3</sup> - Ease of use <sup>4</sup> - Satisfaction <sup>12</sup> - Open comments <sup>10,15</sup> - Overall experience <sup>3,6,7,10,12,14</sup> - Acceptance / preference <sup>4</sup> - Conversational capability <sup>4,6,13,14</sup> - Perceived quality and trust <sup>9</sup> - Suggestions for improvement <sup>1,2,10,15</sup>		- Usability <sup>21</sup> - Feasibility <sup>27</sup> - Usefulness <sup>16</sup> - Ease of use <sup>25</sup> - Satisfaction <sup>17,24-26</sup> - Overall experience <sup>16,17,20,25</sup> - Other open comments <sup>20</sup> - Conversational capability <sup>25</sup> - Perceived quality and trust <sup>20</sup> - Acceptance / preference <sup>20,21,23</sup> - Suggestions for improvement <sup>16,24</sup>		- Usefulness <sup>35</sup> - Ease of use <sup>35</sup> - Satisfaction <sup>30</sup> - Overall experience <sup>28,30,34</sup> - Acceptance / preference <sup>35</sup> - Conversational capability <sup>35</sup> - Suggestions for improvement <sup>30</sup>		- Satisfaction <sup>37</sup> - Overall experience <sup>37</sup>	
	Safety and information quality	- Risk of causing death <sup>12</sup> - Risk of misinformation * - Risk of misunderstanding * - Risk of unintended harms <sup>12</sup>		- Risk of misinformation * - Risk of misunderstanding * - Risk of unintended harms *		-		- CA response capability <sup>36,38,40,43</sup> - Risk of misinformation <sup>40</sup> - Risk of unintended harms <sup>43</sup> - CA response appropriateness <sup>36,38,39,41</sup> - Resources and contents quality <sup>38,43</sup>	
	Functionality	- Response speed <sup>12</sup> - Task achievements <sup>7,12,14</sup> - Engagement functions <sup>8,12,14</sup> - Classification accuracy <sup>1,7</sup> - Understanding and accurate responses <sup>4</sup> - Stability, consistency, and standard *		- Understanding and accurate responses <sup>16</sup>		-		-	