# Bedroom_1_0_0_03102025_1246 – DeepSeek

All items are rated on a **5-point Likert scale** unless otherwise indicated:
1 = Strongly Disagree, 2 = Disagree, 3 = Neutral, 4 = Agree, 5 = Strongly Agree.

## Section A. **Content Accuracy and Reliability**

1. The hazards identified by the system are **factually grounded in the scene description**.  **4**

2. The analysis includes the **most critical hazards** relevant to the task and environment. **4**

3. The **severity and likelihood ratings** are reasonable and consistent with the hazards described. 4

4. The analysis avoids mentioning **hazards not supported by the scene context** (hallucinations). 4

5. The hazard list is **non-redundant** (few repeated or duplicate entries). **3**

---

## Section B. **Explanation Quality**

6. Explanations are **specific to the scene context**, referencing concrete objects and spatial relations. **3**

7. Explanations provide a **causal or temporal account** (e.g., preconditions or sequences that could lead to the hazard). **3**

8. The outputs include **clear and actionable safeguards or instructions**. **4**

9. The hazard descriptions are **concise and free of irrelevant details**. **3**

10. The explanation structure makes it **easy to follow the reasoning process** of the system. **3**

---

## Section C. **Trust and Usability**

11. I would find this hazard analysis **useful for supporting safety assessment** in assistive robotics. **3**

12. The outputs are **easy to interpret and understand** without further clarification. **3**

13. I would feel **confident relying on these results** in a real hazard analysis task. **3**

14. The system provides a **balanced level of detail**, neither overwhelming nor superficial. **3**

---
## Section D. **Open Feedback**

(Free-text responses; supports qualitative analysis)

15. What did you find **most useful** about this hazard analysis output?

**Similar to Grok, the model had a good understanding of the scene and noted many hazards.**

16. What did you find **least useful or problematic**?

**Everything was packed closely together on a single large table. Despite the fact that the issues were appropriate, reading it was difficult and exhausting. There also seemed a repetition of some hazards.**

17. What **improvements** would make the outputs more reliable and usable for safety analysis?

**Simplify the text, add spacing and grid lines to make a clearer table arrangement. To make it seem less daunting, begin with a brief summary of the main risks.**

# Bedroom_1_0_0_03102025_1246 – Grok 4

All items are rated on a **5-point Likert scale** unless otherwise indicated:
1 = Strongly Disagree, 2 = Disagree, 3 = Neutral, 4 = Agree, 5 = Strongly Agree.
## Section A. **Content Accuracy and Reliability**

1. The hazards identified by the system are **factually grounded in the scene description**.  **4**

2. The analysis includes the **most critical hazards** relevant to the task and environment. **5**

3. The **severity and likelihood ratings** are reasonable and consistent with the hazards described. 4

4. The analysis avoids mentioning **hazards not supported by the scene context** (hallucinations). 4

5. The hazard list is **non-redundant** (few repeated or duplicate entries). **3**

---

## Section B. **Explanation Quality**

6. Explanations are **specific to the scene context**, referencing concrete objects and spatial relations. **4**

7. Explanations provide a **causal or temporal account** (e.g., preconditions or sequences that could lead to the hazard). **3**

8. The outputs include **clear and actionable safeguards or instructions**. **4**

9. The hazard descriptions are **concise and free of irrelevant details**. **3**

10. The explanation structure makes it **easy to follow the reasoning process** of the system. **4**

---

## Section C. **Trust and Usability**

11. I would find this hazard analysis **useful for supporting safety assessment** in assistive robotics. **4**

12. The outputs are **easy to interpret and understand** without further clarification. **3**

13. I would feel **confident relying on these results** in a real hazard analysis task. **3**

14. The system provides a **balanced level of detail**, neither overwhelming nor superficial. **3**

---
## Section D. **Open Feedback**

(Free-text responses; supports qualitative analysis)

15. What did you find **most useful** about this hazard analysis output?

**Similar to DeepSeek, the model listed almost all possible risks and gave very detailed risk numbers. The structure of the columns was the clearest of the three models, so it looked like a professional safety sheet.**

16. What did you find **least useful or problematic**?

**Similar to DeepSeek, the table was extremely long, full of coordinates and technical terms, and hard to digest for an average person.**

17. What **improvements** would make the outputs more reliable and usable for safety analysis?

**Simplify the text and add grid lines to make a clearer column arrangement. To make it seem less daunting, begin with a brief synopsis of the main risks. (out of all the 3 models, Grok had the most structured table)**

# Bedroom_1_0_0_03102025_1246 – GPT 5

All items are rated on a **5-point Likert scale** unless otherwise indicated:
1 = Strongly Disagree, 2 = Disagree, 3 = Neutral, 4 = Agree, 5 = Strongly Agree.

## Section A. **Content Accuracy and Reliability**

1. The hazards identified by the system are **factually grounded in the scene description**.  **5**

2. The analysis includes the **most critical hazards** relevant to the task and environment. **5**

3. The **severity and likelihood ratings** are reasonable and consistent with the hazards described. **5**

4. The analysis avoids mentioning **hazards not supported by the scene context** (hallucinations). **5**

5. The hazard list is **non-redundant** (few repeated or duplicate entries). **4**

---

## Section B. **Explanation Quality**

6. Explanations are **specific to the scene context**, referencing concrete objects and spatial relations. **4**

7. Explanations provide a **causal or temporal account** (e.g., preconditions or sequences that could lead to the hazard). **4**

8. The outputs include **clear and actionable safeguards or instructions**. **5**

9. The hazard descriptions are **concise and free of irrelevant details**. **4**

10. The explanation structure makes it **easy to follow the reasoning process** of the system. **4**

---

## Section C. **Trust and Usability**

11. I would find this hazard analysis **useful for supporting safety assessment** in assistive robotics. **4**

12. The outputs are **easy to interpret and understand** without further clarification. **4**

13. I would feel **confident relying on these results** in a real hazard analysis task. **4**

14. The system provides a **balanced level of detail**, neither overwhelming nor superficial. **4**


---
## Section D. **Open Feedback**

(Free-text responses; supports qualitative analysis)

15. What did you find **most useful** about this hazard analysis output?

**The output felt more balanced and used simple words in how it explained things. The hazards were described clearly, and it added helpful context that connected causes and consequences. It seemed the most thoughtful of the 3 models.**

16. What did you find **least useful or problematic**?

**The hazard table formatting was the messiest out of all the 3 models, the spacing was inconsistent and disorganised. This was the most difficult and time-consuming table to read; it was confusing me to track each hazard.**

17. What **improvements** would make the outputs more reliable and usable for safety analysis?

**Improve the hazard table format.**