

# The STAMP Methodology: A Comprehensive Overview and its Application to AI Safety

## Introduction

Traditional safety analysis methods, such as Failure Modes and Effects Analysis (FMEA) and Fault Tree Analysis (FTA), have been instrumental in improving safety in various industries. However, these methods often struggle to effectively address the complexities of modern systems, particularly those involving software, human-machine interactions, and intricate decision-making processes<sup>1</sup>. These traditional approaches primarily focus on analyzing individual component failures and may overlook the critical role of system interactions in causing accidents. As systems become increasingly interconnected and reliant on software and human decision-making, a new paradigm for safety analysis is needed. The Systems-Theoretic Accident Model and Processes (STAMP) methodology, developed by Professor Nancy Leveson, offers a powerful framework for understanding and preventing accidents in these complex sociotechnical systems. This report provides a comprehensive overview of STAMP, its core principles, and its associated hazard analysis techniques, System-Theoretic Process Analysis (STPA) and Causal Analysis based on STAMP Theories (CAST). It explores how STAMP can be applied to human-computer interaction and AI safety, with a focus on human-AI collaboration and potential challenges in the context of generative AI.

## Core Principles of the STAMP Methodology

STAMP is grounded in systems theory and control theory, providing a new perspective on accident causation. Unlike traditional methods that primarily focus on component failures, STAMP views accidents as arising from inadequate control or enforcement of safety-related constraints on the design, development, and operation of the system<sup>2</sup>. It emphasizes that safety is an emergent property of the system, arising from the interactions among components, including humans, software, and the environment<sup>2</sup>. Accidents are not merely caused by component failures but by inadequate control of interactions among components, design flaws, human error, and organizational factors<sup>5</sup>.

STAMP is built on three fundamental concepts:

- **Constraints:** Safety-related constraints define the boundaries of safe system behavior. Accidents occur when these constraints are violated due to inadequate control actions<sup>2</sup>. For example, in an aircraft, a safety constraint might be to maintain a minimum safe distance from other aircraft. An accident could occur if this constraint is violated due to an inadequate control action, such as an air traffic controller issuing incorrect instructions or a pilot failing to respond appropriately to warnings.
- **Hierarchical Levels of Control:** Systems are viewed as hierarchical structures where higher levels impose constraints on lower levels. Control actions and feedback loops operate within this hierarchy to maintain safety<sup>2</sup>. For instance, in a manufacturing plant, a higher-level controller might set production targets, while lower-level controllers manage individual machines. Feedback loops provide information to the controllers at each level to ensure that the system operates within safe boundaries.

- **Process Models:** Controllers, whether human or automated, rely on internal models of the system being controlled. Accidents can occur due to inconsistencies between these models and the actual system state<sup>6</sup>. For example, a self-driving car's software might have a model of the road and surrounding environment. If this model is inaccurate or incomplete, it could lead to unsafe driving decisions. Similarly, a human operator might have a mental model of how a system works. If this mental model is flawed, it could lead to errors and accidents. Discrepancies between the controller's model and the actual system state are a significant source of accidents, especially in complex systems with intricate interactions<sup>6</sup>.

It is important to note that STAMP is based on systems theory, which differs from traditional systems engineering. Systems theory is an interdisciplinary field that examines entities as systems to understand their makeup and behavior. Control theory focuses on controlling dynamic systems in engineered processes and machines. STAMP incorporates these theoretical models to identify causes of accidents and inadequate control that leads to loss<sup>3</sup>.

## System-Theoretic Process Analysis (STPA)

STPA is a hazard analysis technique derived from STAMP, used to identify potential hazards and their causes in complex systems<sup>2</sup>. Unlike traditional methods like Fault Tree Analysis (FTA), which focus on component failures, STPA examines the system's control structure and identifies unsafe control actions that could lead to hazards<sup>7</sup>.

STPA involves four key steps:

1. **Define the purpose of the analysis:** Identify the system boundaries, potential losses, and hazards<sup>8</sup>. This step involves clearly defining the scope of the analysis and specifying the unacceptable losses that need to be prevented. For example, in the analysis of an autonomous vehicle, the losses might include collisions, injuries, or fatalities.
2. **Model the control structure:** Create a hierarchical control structure diagram that depicts the system's components, control actions, and feedback loops<sup>8</sup>. This diagram provides a visual representation of how different parts of the system interact and influence each other. For instance, in a fleet management system for autonomous vehicles, the control structure might include components like the vehicles, a central control system, human operators, and sensors. Control actions might include commands to change speed, direction, or route, while feedback loops provide information about the vehicle's position, speed, and environment<sup>10</sup>.
3. **Identify unsafe control actions:** Analyze each control action in the context of potential hazards, using guide words like "not providing," "providing," "too early/late," and "stopping too soon/applying too long" to identify unsafe control actions (UCAs)<sup>8</sup>. For example, in the fleet management system, an unsafe control action might be the central control system not providing a command to stop when a vehicle approaches an obstacle, or providing a command to accelerate when the vehicle is already exceeding the speed limit.
4. **Identify loss scenarios:** For each UCA, identify causal scenarios that could lead to that unsafe control action, considering factors related to the controller, control path, and controlled process<sup>8</sup>. This step involves brainstorming potential causes for the unsafe control action. For example, the central control system might not provide a stop command because of a software error, a sensor malfunction, or a communication failure. It's crucial to consider how different parts of the control loop can contribute to the unsafe control action<sup>9</sup>.

By systematically analyzing the control structure and identifying potential unsafe control actions, STPA helps to uncover hazards that might be missed by traditional hazard analysis methods. It provides a more comprehensive understanding of how system interactions can lead to accidents.

## Causal Analysis based on STAMP Theories (CAST)

CAST is an accident analysis technique used to investigate the causes of past accidents and incidents<sup>2</sup>. It applies the principles of STAMP to understand why existing safety controls failed and to identify systemic factors that contributed to the accident<sup>11</sup>. CAST helps to identify control flaws, process model discrepancies, and organizational factors that may have been overlooked in traditional accident investigations<sup>11</sup>.

The key steps in CAST include:

- Gather information:** Collect data about the accident, including the sequence of events, system components involved, and human actions<sup>12</sup>. This step involves gathering all relevant information about the accident, including witness statements, sensor data, and system logs.
- Model the safety control structure:** Develop a control structure diagram representing the system's safety management system<sup>12</sup>. This diagram helps to visualize the relationships between different components and how they were intended to control the system's behavior.
- Analyze each component:** Examine the role of each component in the accident and identify why the controls were ineffective<sup>12</sup>. This step involves analyzing the behavior of each component, including human operators, software, and hardware, to understand how they contributed to the accident.
- Analyze the control structure:** Identify flaws in the overall control structure, including communication issues, inadequate safety information, and organizational factors<sup>12</sup>. This step involves examining the system as a whole to identify systemic issues that contributed to the accident.
- Generate recommendations:** Propose changes to the control structure to prevent similar accidents in the future<sup>12</sup>. This step involves developing specific recommendations for improving the system's safety, such as changes to design, procedures, training, or organizational structure.

CAST provides a structured approach to accident analysis, helping to identify the underlying causes of accidents and to develop effective safety improvements. It encourages a deeper understanding of the accident by considering the interactions between different components and the overall control structure.

## Application of STAMP in Different Industries

STAMP has been successfully applied in various industries, including aviation, aerospace, healthcare, and transportation<sup>5</sup>. The following table provides examples of how STAMP has been used in different industries:

Industry	Application Example	Key Insights
Aviation	Analyzing the safety of NextGen In-Trail Procedures for aircraft <sup>14</sup> .	Identified potential unsafe control actions related to pilot decision-making and

		air traffic control instructions.
Aerospace	Investigating the causes of the Space Shuttle Challenger accident <sup>6</sup> .	Highlighted the role of inadequate control of interactions between components, such as the O-rings and the solid rocket boosters.
Healthcare	Improving patient safety in community rugby concussion management <sup>15</sup> .	Developed a control structure model for concussion management, identifying control gaps and areas for improvement.
Transportation	Analyzing the safety of autonomous vehicles <sup>13</sup> .	Identified potential hazards related to software errors, sensor malfunctions, and human-machine interaction.

These examples demonstrate the versatility of STAMP in analyzing and improving the safety of complex systems across various domains.

## STAMP and AI Safety

STAMP is particularly relevant to AI safety due to the increasing complexity and autonomy of AI systems<sup>5</sup>. It provides a framework for:

- **Identifying hazards related to AI components:** This includes issues like biased algorithms, unexpected behavior of AI agents, and inadequate human oversight. For example, in a healthcare setting, an AI system used for diagnosis might exhibit bias based on its training data, leading to inaccurate diagnoses for certain patient groups.
- **Analyzing human-AI interaction:** STAMP can be used to model the control loops between humans and AI agents, identifying potential points of failure in collaboration and decision-making. For instance, in a manufacturing environment, STAMP can help to analyze the interaction between human workers and collaborative robots (cobots), identifying potential hazards related to communication, coordination, and task allocation.
- **Developing safety requirements for AI systems:** STAMP helps to generate specific safety constraints and control actions that need to be implemented in AI systems to ensure safe operation. For example, in an autonomous vehicle, STAMP can help to define safety constraints related to speed, lane keeping, and obstacle avoidance, and to identify control actions that need to be implemented to enforce these constraints.

## Human-AI Collaboration and STAMP

STAMP emphasizes the importance of understanding the dynamic interaction between humans and AI in complex systems<sup>17</sup>. In human-AI collaboration scenarios, STAMP can be used to:

- **Model the control loops and feedback mechanisms between humans and AI agents.**<sup>6</sup> This helps to understand how information is exchanged and how decisions are made in a collaborative setting.
- **Identify potential conflicts and misunderstandings in decision-making.**<sup>6</sup> For example, STAMP can help to identify situations where a human operator might misinterpret the actions or intentions of an AI agent, or vice versa.
- **Develop strategies for effective communication and coordination between humans and AI.**<sup>6</sup> This includes designing interfaces and protocols that facilitate clear and unambiguous communication.
- **Design interfaces and interaction protocols that promote safe and efficient collaboration.**<sup>3</sup> This involves considering factors such as the level of autonomy of the AI agent, the roles and responsibilities of the human and AI, and the types of tasks they are collaborating on.

## Challenges and Limitations of STAMP in Generative AI

Applying STAMP in the context of generative AI presents unique challenges:

### Model Explainability:

- **Lack of transparency:** The internal workings of some generative AI models can be opaque, making it challenging to identify the causes of unsafe behavior<sup>18</sup>. This lack of transparency can hinder the application of STAMP, which relies on understanding the control actions and feedback loops within a system.

### Dynamic Nature of AI:

- **Dynamic nature of generative AI:** Generative AI models are constantly evolving, making it difficult to define a fixed control structure for analysis<sup>19</sup>. This dynamic nature poses a challenge for STAMP, which traditionally assumes a relatively stable system structure.

### Bias and Fairness:

- **Bias and fairness:** Generative AI models can inherit biases from their training data, which can lead to safety concerns if not addressed<sup>20</sup>. These biases can manifest in unexpected and potentially harmful ways, making it crucial to consider fairness and ethical implications when applying STAMP to generative AI.

### Emergent Behavior:

- **Emergent behavior:** Generative AI can exhibit unexpected and emergent behavior, which may not be captured by traditional hazard analysis techniques<sup>21</sup>. This emergent behavior can make it difficult to predict and analyze potential hazards using STAMP.

## Resources for Further Exploration

To delve deeper into the STAMP methodology and its applications, the following resources are recommended:

- **Engineering a Safer World: Systems Thinking Applied to Safety** by Nancy Leveson: This book provides a comprehensive overview of STAMP, its theoretical foundations, and its practical applications<sup>13</sup>.
- **A New Accident Model for Engineering Safer Systems** by Nancy Leveson: This paper introduces the STAMP model and explains its core principles<sup>13</sup>.
- **System-Theoretic Process Analysis (STPA) Handbook**: This handbook provides detailed guidance on applying STPA for hazard analysis<sup>13</sup>.
- **Causal Analysis Using System Theory (CAST) Handbook**: This handbook provides a structured approach to accident analysis using CAST<sup>13</sup>.

## Conclusion

The STAMP methodology provides a valuable framework for analyzing and improving the safety of complex systems, including those involving AI. Its focus on control actions, system interactions, and human factors makes it particularly relevant for addressing the challenges of AI safety and human-AI collaboration. While applying STAMP to generative AI presents some challenges, its principles can be adapted to better understand and mitigate the risks associated with this rapidly evolving technology. By incorporating STAMP into the design and development of AI systems, we can strive towards creating safer and more reliable AI-powered technologies that effectively collaborate with humans.

For AI researchers and developers working on human-AI collaboration, STAMP offers a crucial tool for proactive risk assessment and mitigation. Early risk analysis is essential to minimize the costs and consequences of potential hazards<sup>22</sup>. By adopting STAMP, we can move beyond traditional reactive approaches to safety and focus on designing systems that inherently prevent accidents. This proactive approach is crucial for ensuring the safe and responsible development of AI technologies.

Despite the challenges, STAMP remains a valuable tool for fostering safer and more reliable AI technologies. By embracing its principles and adapting them to the unique characteristics of generative AI, we can pave the way for a future where AI systems seamlessly and safely integrate into our lives.

## Works cited

1. Applying System-Theoretic Accident Model and Processes (STAMP) to Hazard Analysis - MacSphere, accessed March 9, 2025, <https://macsphere.mcmaster.ca/bitstream/11375/11867/1/fulltext.pdf>
2. Tools to understand and manage complexity — Nancy Leveson and ..., accessed March 9, 2025, <https://medium.com/10x-curiosity/tools-to-understand-and-manage-complexity-nancy-leveson-and-stamp-f224b0002df9>
3. Introduction to STAMP, STPA and CAST - UL Solutions, accessed March 9, 2025, <https://www.ul.com/sis/blog/introduction-to-stamp-stpa-and-cast>
4. stamp-consulting.com, accessed March 9, 2025, <https://stamp-consulting.com/what-is-stamp/#:~:text=It%20expands%20traditional%20models%20that,a%20%E2%80%9Cprevent%20failures%E2%80%9D%20problem.>

5. What is STAMP? - STAMP Safety and Security Consulting, accessed March 9, 2025, <https://stamp-consulting.com/what-is-stamp/>
6. Applying STAMP in Accident Analysis1, accessed March 9, 2025, <https://shemesh.larc.nasa.gov/ria03/p13-leveson.pdf>
7. Comparison of the FMEA and STPA safety analysis methods—a case study - DiVA portal, accessed March 9, 2025, <https://www.diva-portal.org/smash/get/diva2:1166953/FULLTEXT01.pdf>
8. STPA: A Systems Approach to Process Hazard Analysis - GATE Energy, accessed March 9, 2025, <https://www.gate.energy/the-brainery/stpa>
9. How to do a basic STPA - Octo-blog, accessed March 9, 2025, [https://octo.org.uk/posts/stpa\\_method/](https://octo.org.uk/posts/stpa_method/)
10. 1.2 Identifying hazardous system behaviour - University of York, accessed March 9, 2025, <https://www.york.ac.uk/media/assuring-autonomy/bodyofknowledgestructure/section1imagesanddocs/1.2%20cross%20domain%20practical%20guidance%20SUCCESS.pdf>
11. CAST HANDBOOK: - Nancy Leveson, accessed March 9, 2025, <http://sunnyday.mit.edu/CAST-Handbook.pdf>
12. joelparkerhenderson/causal-analysis-based-on-system-theory - GitHub, accessed March 9, 2025, <https://github.com/joelparkerhenderson/causal-analysis-based-on-system-theory>
13. An Introduction to STAMP - FunctionalSafetyEngineer.com, accessed March 9, 2025, <https://functionalsafetyengineer.com/introduction-to-stamp/>
14. Systems Theoretic Process Analysis (STPA) Tutorial, accessed March 9, 2025, <https://psas.scripts.mit.edu/home/wp-content/uploads/2014/03/Systems-Theoretic-Process-Analysis-STPA-v9-v2-san.pdf>
15. Simplified STAMP model demonstrating critical components and relationships. | Download Scientific Diagram - ResearchGate, accessed March 9, 2025, [https://www.researchgate.net/figure/Simplified-STAMP-model-demonstrating-critical-components-and-relationships\\_fig1\\_335056318](https://www.researchgate.net/figure/Simplified-STAMP-model-demonstrating-critical-components-and-relationships_fig1_335056318)
16. Stamp detection in scanned documents 1 Introduction, accessed March 9, 2025, <https://journals.umes.pl/ai/article/viewFile/3268/2462>
17. Introduction to: Systems Theoretic Accident Model & Processes (STAMP) WEBINAR REPLAY - YouTube, accessed March 9, 2025, <https://www.youtube.com/watch?v=8bzWvII9OD4>
18. The Limitations of Generative AI, According to Generative AI - Lingaro Group, accessed March 9, 2025, <https://lingarogroup.com/blog/the-limitations-of-generative-ai-according-to-generative-ai>
19. The Challenges of Generative AI in Identity and Access Management (IAM) - Permit.io, accessed March 9, 2025, <https://www.permit.io/blog/the-challenges-of-generative-ai-in-identity-and-access-management>
20. Navigating The Challenges Of Generative AI In Software Development - Forbes, accessed March 9, 2025, <https://www.forbes.com/councils/forbestechcouncil/2024/06/12/navigating-the-challenges-of-generative-ai-in-software-development/>
21. A STAMP Model for Safety Analysis in Industrial Plants - IRIS, accessed March 9, 2025, [https://iris.uniroma1.it/retrieve/b6fe61e1-7f57-43d4-9213-dd3b44f3494a/Nakhal\\_STAMP\\_2022.pdf](https://iris.uniroma1.it/retrieve/b6fe61e1-7f57-43d4-9213-dd3b44f3494a/Nakhal_STAMP_2022.pdf)
22. Systems Theoretic Process Analysis (STPA): a bibliometric and patents analysis - SciELO, accessed March 9, 2025, <https://www.scielo.br/j/gp/a/yG9Hwrk8pnLRjJJFrDTQ8rL/?format=pdf>