5.3  Safety & Risk Reduction:
 With critical applications like self driving vehicles, financial systems, and legal decision-making, unsafe AI behaviour can create "catastrophic consequences".
Reasoning frameworks:
• Help the AI avoid risky paths or other shortcut.
• Logically consistent in an presented situations even with edge cases.
• Provide an avenue for AI to align with ethical norms through human moral reasoning
5.4 Enabling Generalization and Flexibility:
 AI models with reasoning are better

suited to follow through with new tasks during operationalization that are not included in their training or even not near that training. Instead of remembering past examples, reasoning ought to: Communication is about cognition and cognition is about:

• Reasoning is about - identify relevant similarities.
• Apply knowledge across different domains (analogical reasoning).
• produce novel solutions in new environments.

## 5.5. Supporting explainable AI (XAI) objectives:

Explainable AI (XAI) is an emerging research domain dedicated to increasing the Interpretability/Understanding of AI systems.

Reasoning is the basis of Explainable AI since:

• users can be shown logical steps
• unreasoned errors can be diagnosed or corrected
• The model can provide justifications, in human terms, for its decisions.

A reasoning algorithm or model will allow for explainability rather than superficial or incomplete explanations.

### 5.6: Ethical and Legal Accountability :

Many regulators, globally (i.e. the EU's AI Act), are requiring AI systems to provide explanation, particularly in high-risk sectors.

Reasoning processes would also allow an AI to:

• explain how it makes decisions - provide an evidence trail - evidence fairness/non-discrimination Protecting stakeholders legally and ethically.

### Summary:

Reasoning transforms AI from a reactive device to a trusted and disciplined decision-maker which can resolve many issues associated with trustworthiness and safety in complex human contexts.

## 6. AI SAFTEY & TRUSTWORTHY:

Artificial Intelligence (AI) safety and trustworthiness are closely related concepts that aim to ensure AI systems are beneficial, reliable, and aligned with human values and societal well-being. Here's a breakdown of what each entails:

### 6.1 AI SAFTEY:

AI safety is a field that focuses on investigating ways to ensure AI systems do not cause unintended harm or act in ways misaligned from human goals, especially as systems

advance and become more and more autonomous. It consists of activities and research geared to limiting risks associated with AI development and deployment. Four of the aspects of AI safety include:

Key Aspects of AI Safety Include:

| Aspect | Explanation |
| --- | --- |
| Robustness | The AI should handle imperfect, adversarial, or unexpected inputs without catastrophic failures. |

| Aspect | Explanation |
| --- | --- |
| Alignment | The AI's goals and behaviors must match human ethical principles, laws, and societal expectations. |
| Controllability | Humans must retain the ability to intervene, correct, or shut down AI systems when needed. |
| Preventing Reward Hacking | AI should not "game" its reward |

signals in ways that cause unintended consequences.

## 6.2 Trustworthy AI:

Trustworthy AI is a broader concept that encompasses AI safety but also includes other crucial dimensions that build confidence in AI systems among users and society. The European Union's Ethics Guidelines for Trustworthy AI, for example, outlines three main components: Core Principles of Trustworthy AI:

Principle Explanation

Transparency
AI decisions should be understandable; no "black box" behavior.

Fairness
AI should not discriminate based on race, gender, or other biases.

Accountability
Developers and users should be able to explain and justify AI behavior.

Privacy and Security
AI must protect

user data and not
expose sensitive
information.
Reliability
AI should perform
accurately and
consistently over
time and across
different
conditions.
Ethical
Alignment
AI behavior should
conform to societal
and moral norms.

### 6.2.1 Why Are AI Safety and Trustworthiness Critical Together?

As AI increasingly influences:

• Healthcare decisions

• Financial approvals

• Hiring and education

• Autonomous vehicles

• Military defense

Mistakes are not just technical errors — they become real-world ethical, legal, and human rights issues.

Thus:

• AI Safety ensures systems do not fail catastrophically.

• Trustworthy AI ensures systems earn and deserve human trust.

Together, they build an AI ecosystem that is beneficial, reliable, and

ethically sound for society.

Summary:

AI Safety makes AI reliable; Trustworthy AI makes it ethically acceptable. Both are essential for a future where AI empowers humanity without harming it.

## 7. Reasoning Contributes to AI Safety and Trustworthiness

### 7.1 Reasoning Increases Transparency:

When AI systems reason explicitly (e.g. demonstrating its logical steps, like a human would), they are more transparent to people.

Humans can:

• Observe how and AI came to their decision

• Understand why they formed certain conclusions

• Retrospectively trace errors in reasoning if something went wrong

Transparency is a critical element of trustworthy AI - without reasoning, AI could be considered a black box, whereby it is unobtainable for people to trust.

### 7.2. Reasoning Improves Robustness and Reliability

Reasoning enables AI to check its intermediate steps, and validate its assumptions and reason logically when new situations arise as opposed to guessing.

This increases:

- Robustness (AI will work even when new, unexpected inputs are introduced)
- Accuracy (AI will generate fewer silly errors)
- Consistency (AI is more likely to do something in the same way over time)

For the purposes of AI safety, robustness and reliability is the key element - being imperfect in high risk domains such as healthcare, aviation and autonomous driving.

3. Reasoning Aids in Error Detection and Error Correction:

If an AI can reason its way through a problem step-wise, it may enable humans (or the AI) to:

- Identify errors early (i.e. incorrect assumptions).
- Understand faulty logic.
- Make corrections to the output before harm is caused.

Overall, this reasoning aids in the self-correction of AI and increases safety, and it aids humans in auditing the behaviour of the AI in an error situation.

7.4 Reasoning Enhances Ethical and Fair Decision:-

Making As an advisor to humans, AI with structured reasoning can:

- Weigh ethical trade-offs of (e.g.. fairness vs accuracy)
- Make a reasonable application of social norms (e.g.,

nondiscrimination)

• Account for the long term effects of its actions

This is what makes AI trustworthy, as it acts in much the way humans' moral beliefs do, instead of simply maximizing for algorithms that prioritize profit or efficiencies.

7.5. Reasoning can reduce bias and random behaviour:

Without reasoning, AI can take "short cuts" by using potentially misleading correlations (E.g., relating gender or race to rank and file performance).

Using structured reasoning AI can:

• Make behavioural choices based upon logical causative and effect instead of just correlations

• Explain logically why an outcome was fair

This implies that AI may have less bias, that it is more fair and as such should be considered more trustworthy and safe for society.

Summary:

Reasoning is a "thinking brain" inside AI -- reasoning limits random guessing, forces logical and ethical limits, improves transparency, and provides a basis to make AI systems safer, fairer and trustworthy for our society.

Contribu

tion

AI

Safety

Benefi

t

Trustworth

iness

Benefit

Transpare

ncy

Easier to

detect

errors

and

Builds user

trust

Contribu

tion

AI

Safety

Benefi

t

Trustworth

iness

Benefit

prevent

harm

Robustnes

s

Handles

novel,

adversar

ial cases

safely

Consistent

behavior

builds

reliability

Error
Detection

Safer
correcti
ons
during
operatio
ns

Auditable and
explainable
decisions

Ethical
Reasoning

Avoids
harmful
or unfair
outcome
s

Aligns with
human values
and ethics

Bias
Reduction

Prevents
catastro
phic
unfairne
ss

Ensures
fairness and
social
acceptance