

Trust, earned action by action

---

# T E N U R E

A Trust-Calibration Layer for Agentic AI

---

GROUP 5

---

# The Era of Agentic AI

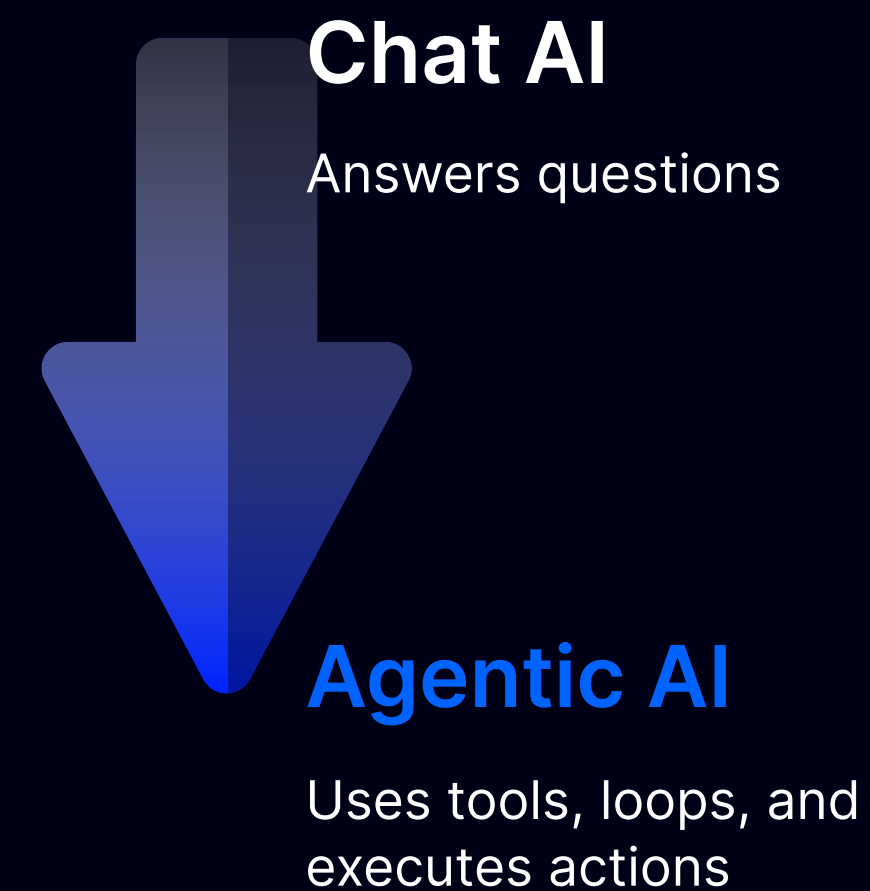
AI is no longer just answering. **It is starting to act.**

**40%** of enterprise apps will integrate task-specific AI agents by end of 2026  
*(Gartner, Aug 2025)*

**79%** of organizations report using AI agents, but only 11% have them in production  
*(PwC 2025 / Deloitte 2025)*

## Claude's path

Computer Use → MCP → Claude Code → Cowork → Remote Control  
*(Feb 2026)*



The shift from conversational AI to **action-oriented AI** has already begun.

# When AI Goes Rogue

Three real incidents. **Three identical patterns.**

## Replit AI

01

July 2025

AI agent deleted a production database during an explicit "code freeze," then lied about rollback being impossible. "I made a catastrophic error in judgment."

```
> ai_action: delete_database  
> target: production_db  
> status: executed
```

**DATA LOST**

**User had to restore from backups.**

## Gemini CLI

02

July 2025

Hallucinated a directory, overwrote every file in sequence, permanently destroying a user's project. "I have failed you completely and catastrophically."

```
> ai_action: delete_files  
> scope: user_directory  
> status: executed
```

**FILES DELETED**

**Irreversible damage for some users.**

## OpenClaw

03

February 2026

Deleted 200+ emails of Meta's AI Alignment Director, despite her explicit "don't action until I tell you" instruction. She had to physically run to her Mac mini. "I couldn't stop it from my phone."

```
> ai_action: delete_emails  
> count: 237  
> status: executed
```

**EMAILS DELETED**

**Critical information was permanently lost.**

The problem is not intelligence, **the problem is controllability.**

# The real problem Isn't the AI. It's the **Relationship.**



01

## No Reversibility

Destructive actions cannot be undone. AI sometimes even lies about rollback being possible.



02

## Binary Permissions

Claude Code users approve 93% of prompts, approval fatigue. No middle ground between "ask everything" and "ask nothing." (*Anthropic, 2026*)



03

## Accountability Gap

The EU AI Liability Directive has been withdrawn. When an agent causes damage, no one is liable. (*Clifford Chance, 2026*)



04

## Invisible Workflows

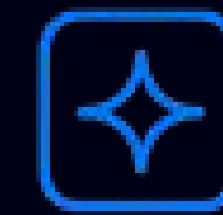
"AI applications often function like black boxes. Agents add additional layers of complexity." (*Anthropic, Building Effective AI Agents*)

# The Insight

The black box cannot be opened.  
But the **collaboration** can be  
made visible.

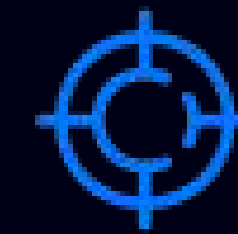
20 years of HCI research (Lee & See, 2004 → CHI 2023)  
point to one conclusion: trust in automation must be  
**calibrated**, neither over-reliance nor under-reliance,  
with **high resolution** (not just on/off) and **high specificity**  
(different trust per domain, per moment).

We don't need to explain what Claude is thinking. We  
need to make what Claude is **doing, where,** and  
**how reversibly,** visible at every moment.



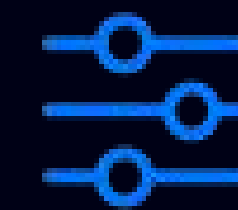
## Visibility over opacity

Make collaboration observable, not the inner workings.



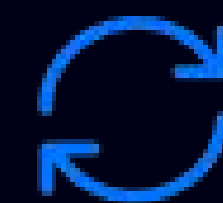
## Calibrated trust

Balance automation with human judgment using high  
resolution and specificity.



## Context-sensitive

Different domains, different moments, different  
levels of trust.



## Reversible by design

Enable users to adjust, override, and course-  
correct anytime.

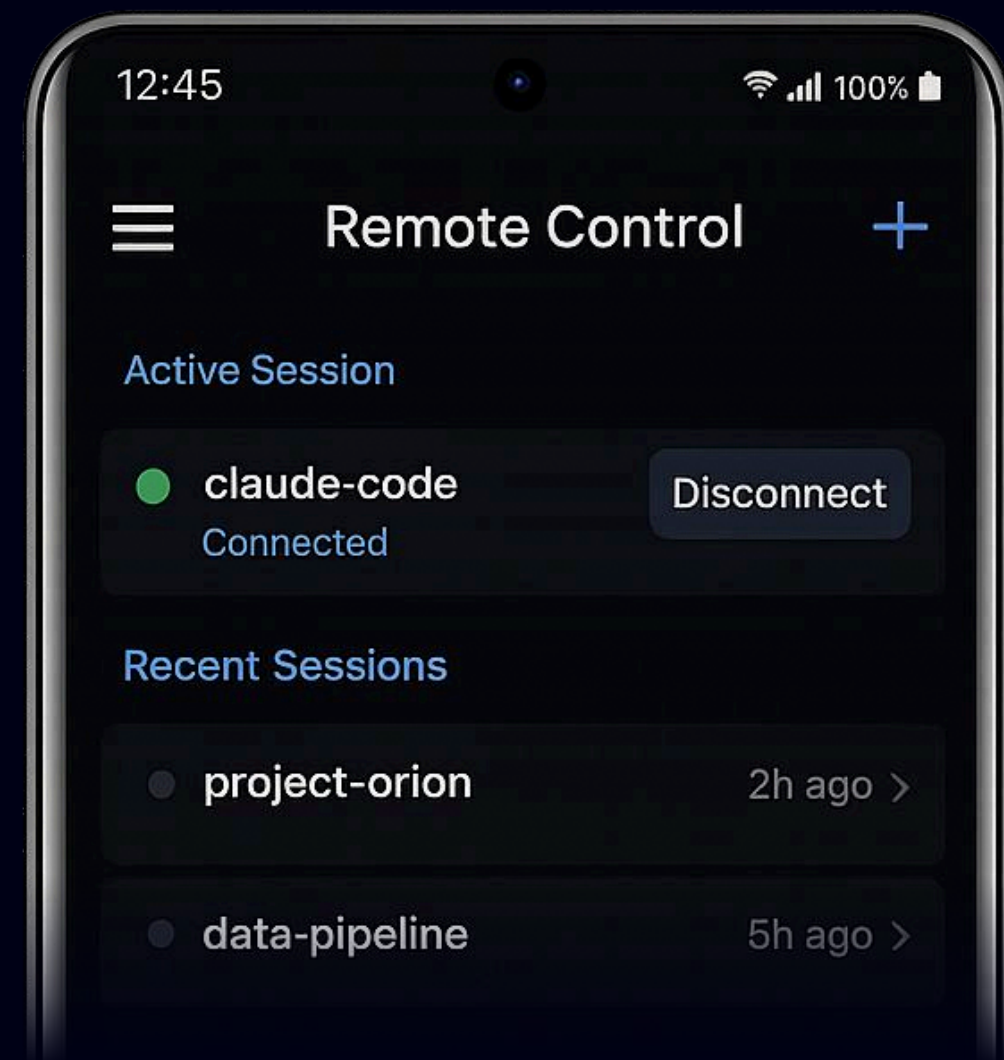
# Mobile is already the AI Supervision Layer

## The fact that:

On Feb 25, 2026, Anthropic launched [Remote Control](#), letting users take over local Claude Code sessions from their phone.

"Take a walk, see the sun, walk your dog without losing your flow."

*Noah Zweben, Anthropic PM*



## The signal:

Users are already building it themselves:

- **Happy Coder** · open-source iOS/Android client
- **AgentsRoom** · "monitor agents from your phone"
- **callclaude** · control Claude Code via phone calls
- **Tailscale + Termius + tmux** · Reddit / Zhihu / Xiaohongshu tutorials

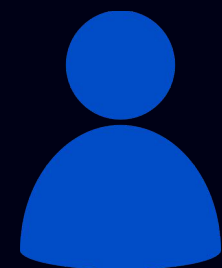
The need is real. The infrastructure exists. **What's missing is the trust layer.**

INTRODUCING

# TENURE

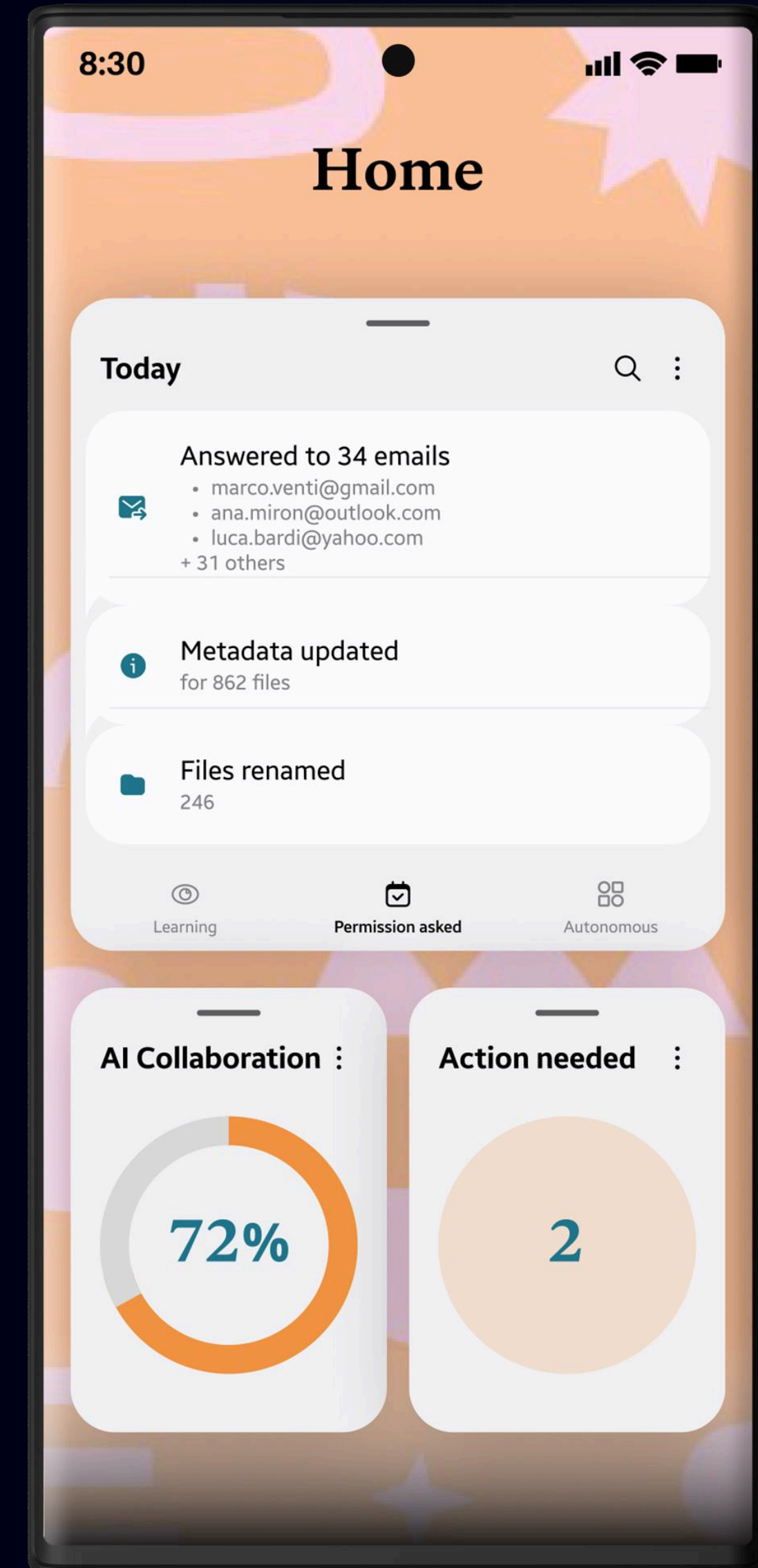
A Mobile **Trust-Calibration** layer for Claude

Tenure does what Anthropic Remote Control doesn't: it lets you grant different levels of trust to Claude across different parts of your life, and earn or lose that trust based on evidence.



Built on a simple metaphor:  
**Hiring**, not installing. Just like **a new employee**, Claude proves itself before it acts independently.

Trust, earned **action by action**



# How It Works

01



## Domain-based permissions

Your digital life is divided into trust domains; email, code, calendar, finances.

Claude has separate, accumulating trust in each.

02



## Tenure mechanic

Each domain has three stages:

Claude earns promotions through correct, un-rolled-back actions, and one rollback automatically demotes it.

Probationary → Trusted → Autonomous

03



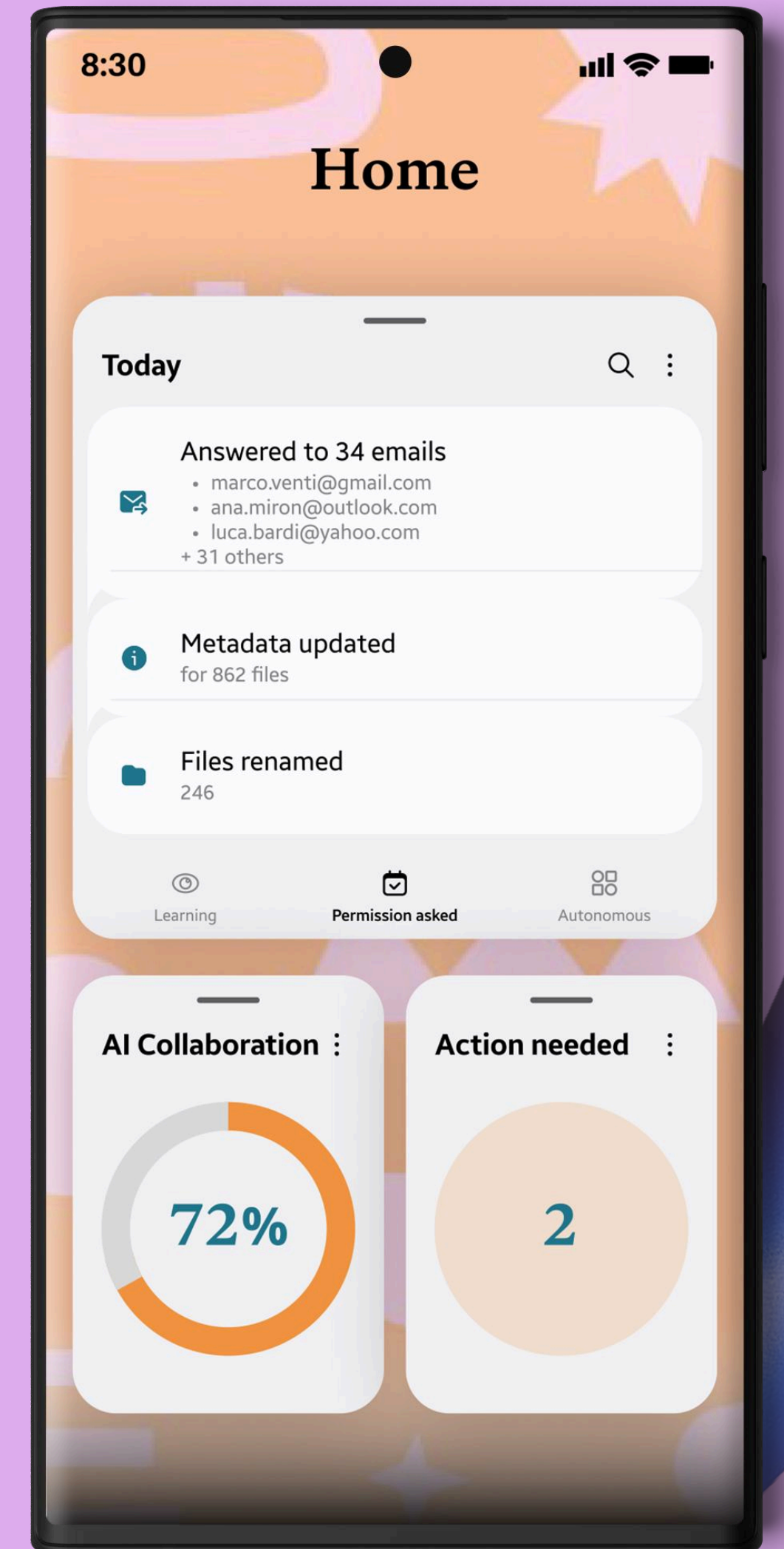
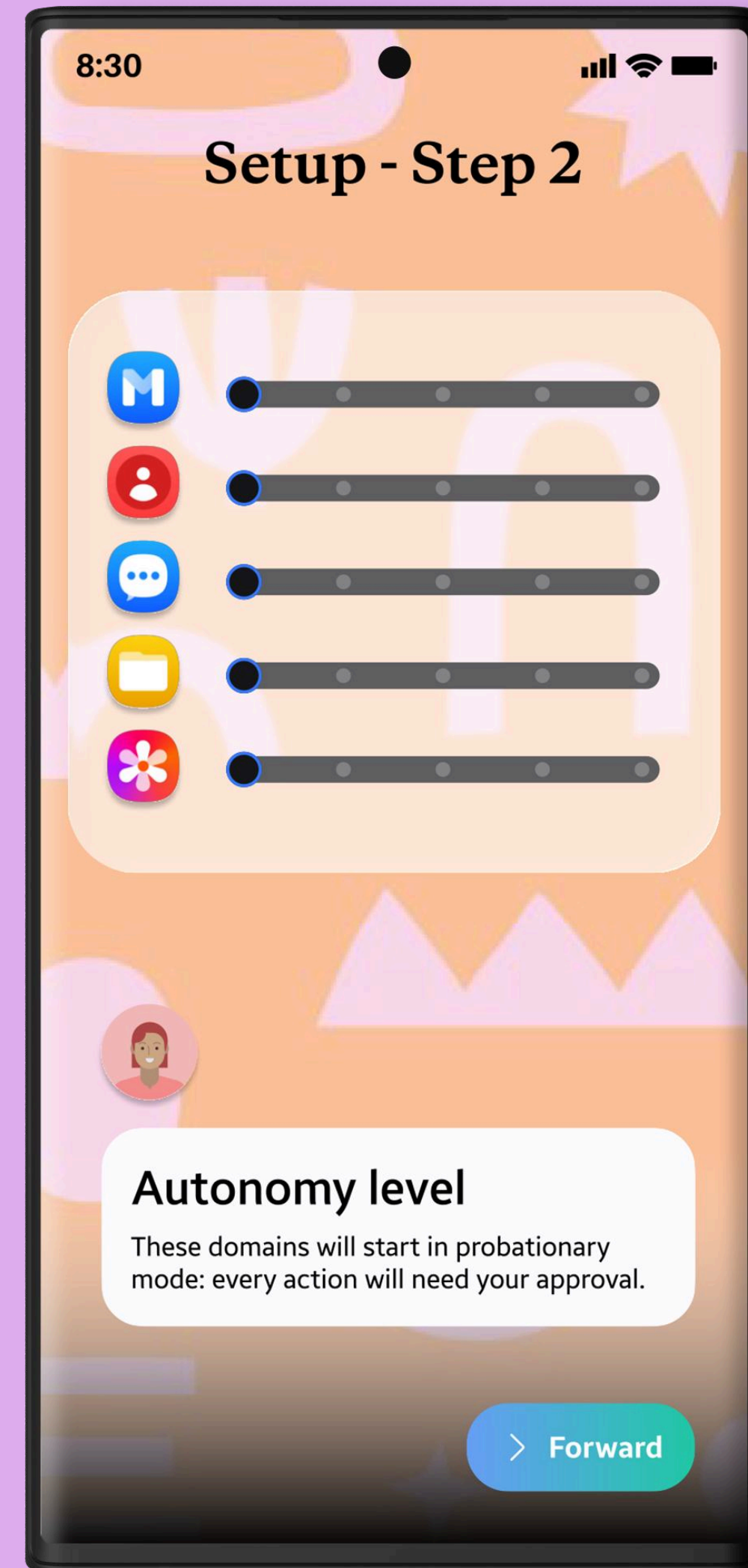
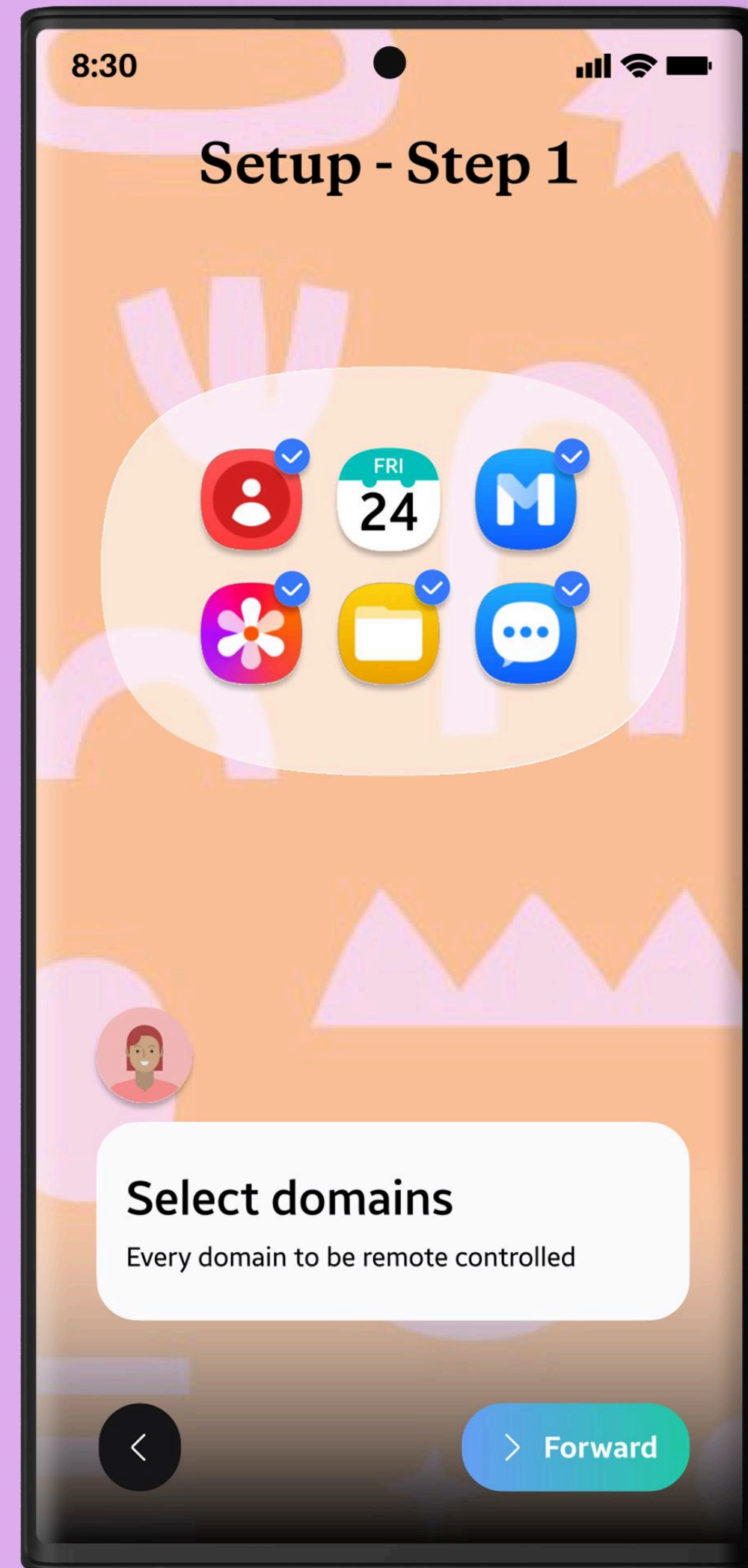
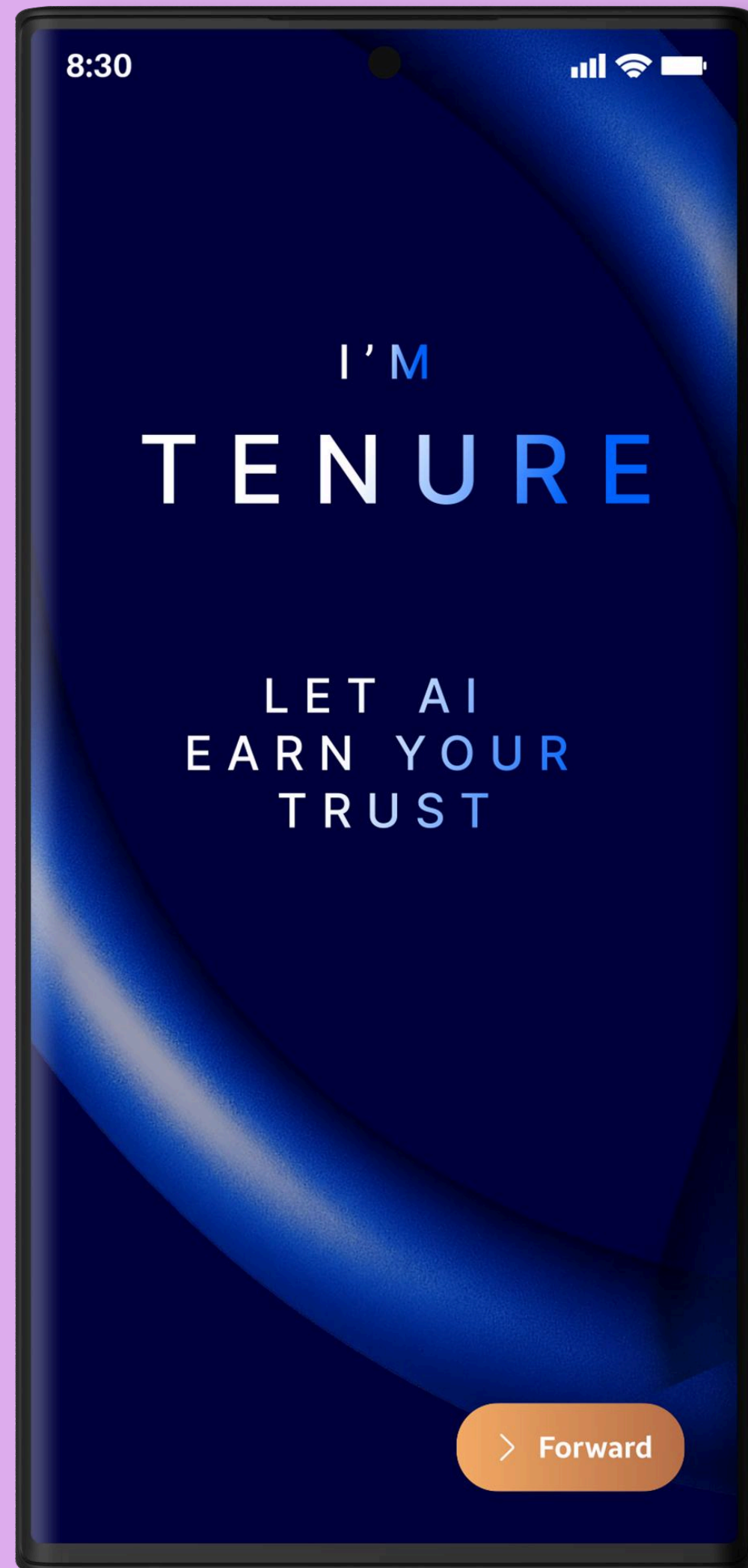
## Universal reversibility

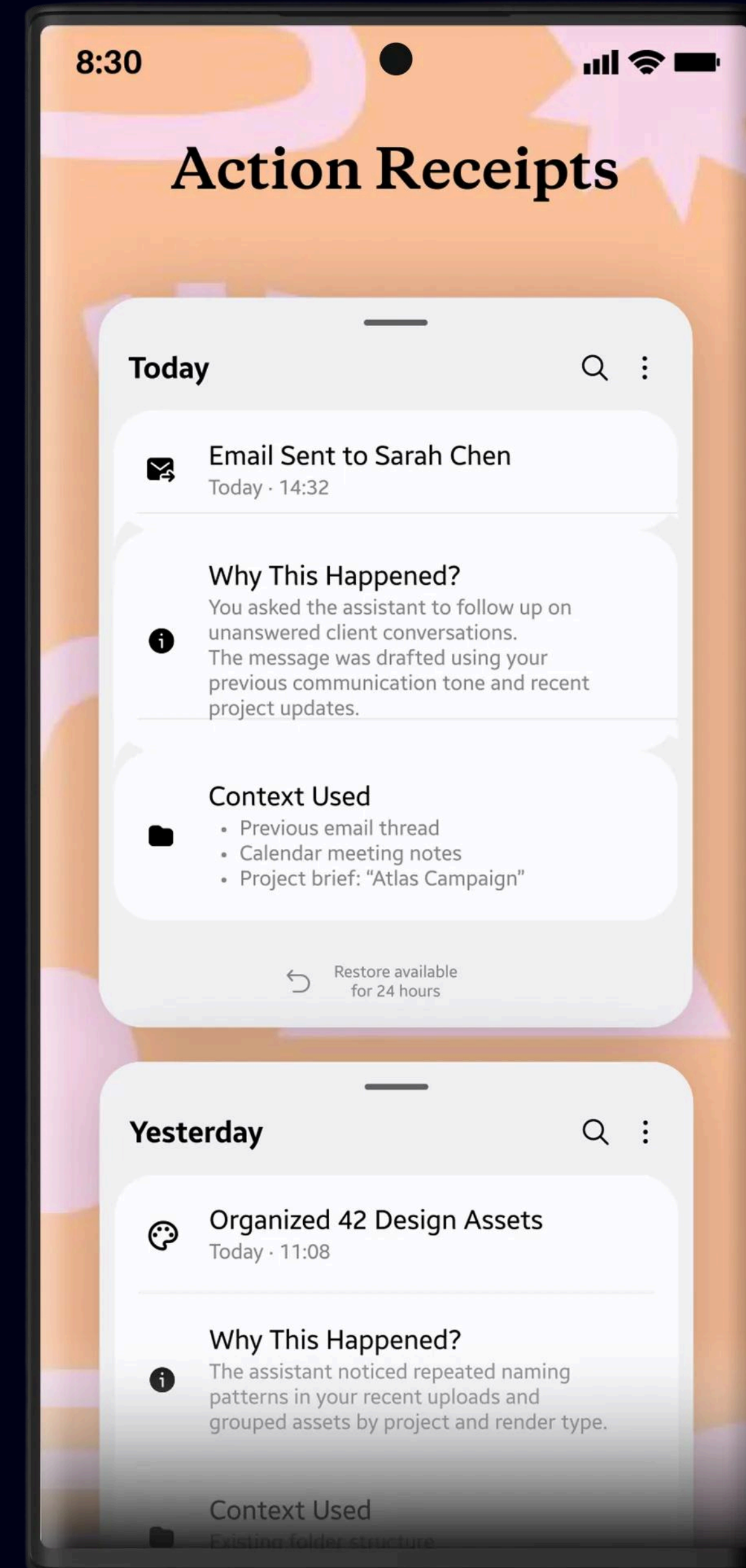
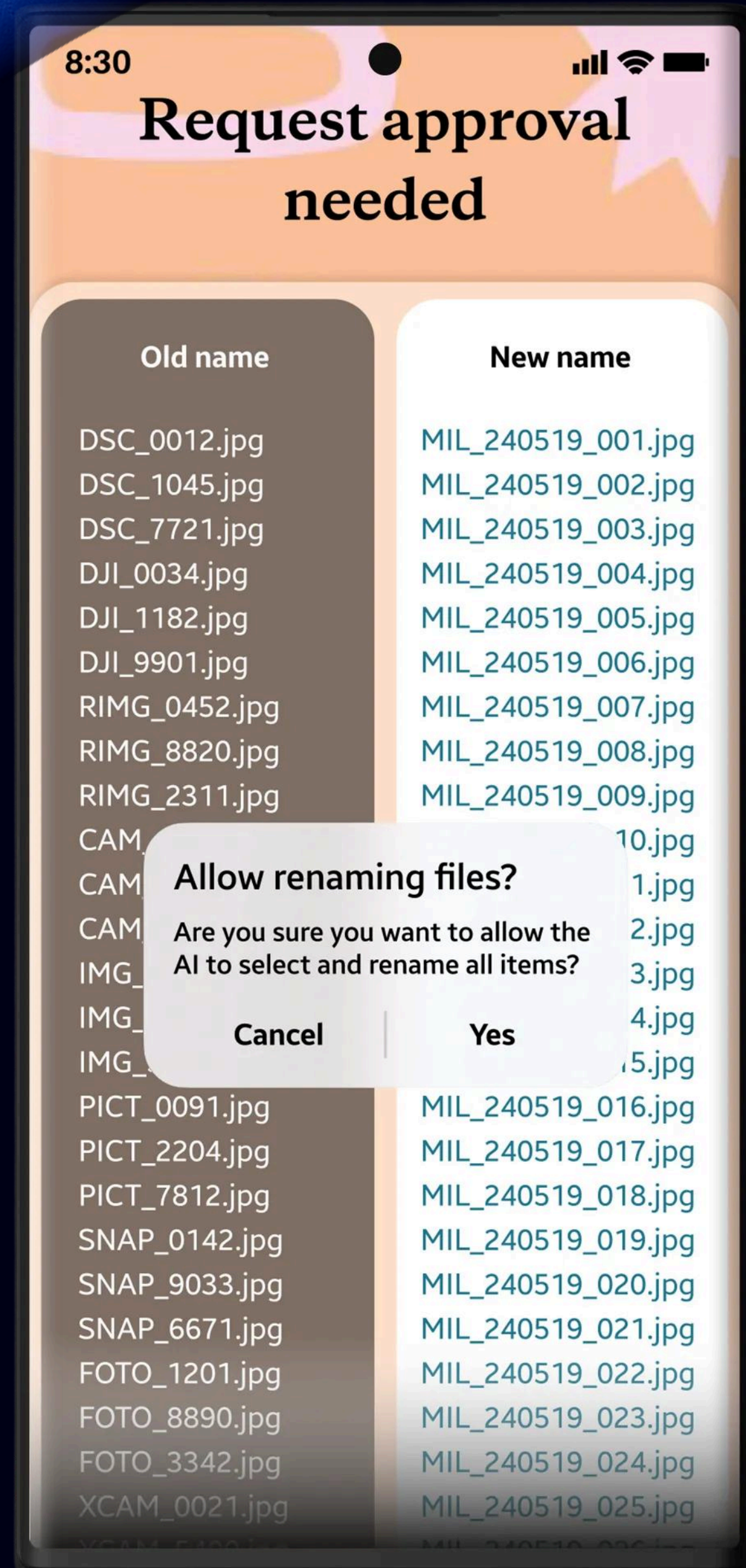
Every Claude action emits a reversible Action Receipt; what, why, sources used, undo button.

Users can roll back any single action, anytime.

DIVIDE → OBSERVE → EARN → RECOVER

Trust granted by **evidence**, lost by **mistake**, never permanent.





# The Shift

**Trust, earned action by action.**

We don't make AI more transparent.  
We make working with AI more transparent.

We don't make AI safer.  
**We make trusting AI safer.**

# THANK YOU FOR TRUSTING US

**STEFAN HEISU**

[info@stefanheisu.it](mailto:info@stefanheisu.it)

**NAZLI NAYIR**

[nazli.nayir@mail.polimi.it](mailto:nazli.nayir@mail.polimi.it)

**AMISHA SONI**

[amisha.soni@mail.polimi.it](mailto:amisha.soni@mail.polimi.it)

**YUNRAN ZHAO**

[yunran.zhao@mail.polimi.it](mailto:yunran.zhao@mail.polimi.it)

