



Autonomous **Trustworthy Agents** (ATA)

06.11.2025

Engineering Kiosk Alps



Otera

David Peer

david.peer@otera.ai

Machine Learning Researcher

“... Zeit für 800 Folien”



Otera

Previously DeepOpinion

Autonomy, Control, and Velocity for processes where quality, consistency, and trust are non-negotiable.

Insurance - Claims Processing



Claims

Alice and her brother Bob booked a trip to Spain. Unfortunately, Bob got the flu.

I lost my passport and, therefore, I was not able to travel.

...

Terms & Conditions (T&C)

Contract

Definition of Words:

Close Relative: Your spouse, partner, parent, parent-in-law, grandparent, child, grandchild, brother, ...

What is covered?

1. The policyholder is covered for any claim arising from or attributable to sickness or death of you, your travelling companion or a close relative...

What is not covered?

1. If you lost your passport and you were not able to go on the trip, you are not covered...

Coverage





“Raw” GenAI for Claims Processing

Claims Processing with "raw" GenAI

Analysis of Coverage

Relative

Mother, mother-in-law, father, father-in-law, sister, sister-in-law, brother, brother-in-law, wife, husband, civil partner, daughter, daughter-in-law, son, son-in-law, grandparent, grandparent-in-law, grandchild, step parent, step child, step sister, step brother, foster child, partner or fiancé/fiancée of an insured person.

the death, bodily injury or illness of:

relation due to the illness of a relative (cousin) and a friend who is covered under Section 1 - whom you are travelling (friend) is covered under Section 1 - no relation to suggest any major exclusion applies.

yes

1. **Cousin got the flu:** A cousin is explicitly defined as a **relative**. Illness of a relative is a covered reason for cancellation. The term "flu" falls under "illness" or "disease".
2. **Close relative and friend got sick:** This would cover the *relative* (cousin) and a *whom you are travelling* (friend).



Enter a prompt for Gemini

+ Tools

Can we really trust¹ "raw" GenAI?





Trustworthy¹ AI

**Accuracy &
Stability**

**Transparency &
Explainability**

**Human Agency &
Oversight**

**Privacy &
Security**

Accountability

**Fairness &
Non-Discrimination**



~~Solution~~ - Generic Scaling Law for Everything

More Data



More Compute



Autonomous Trustworthy Agents (**ATA**)

Idea...

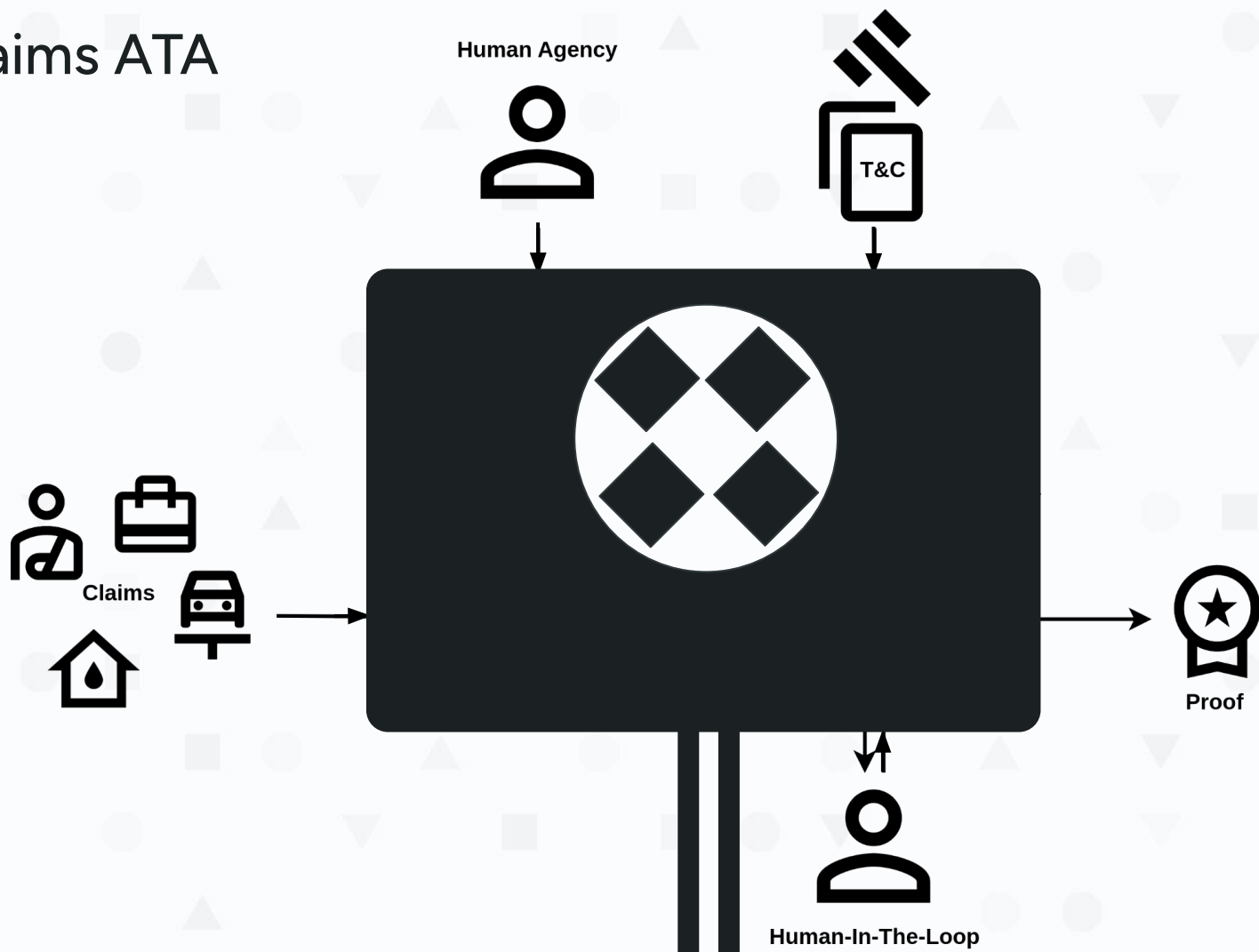
Maybe we can use **Many-Sorted First-Order Logic** to formalize T&Cs?

Knowledge Base \rightarrow T&C, Claim
Goal \rightarrow Covered?

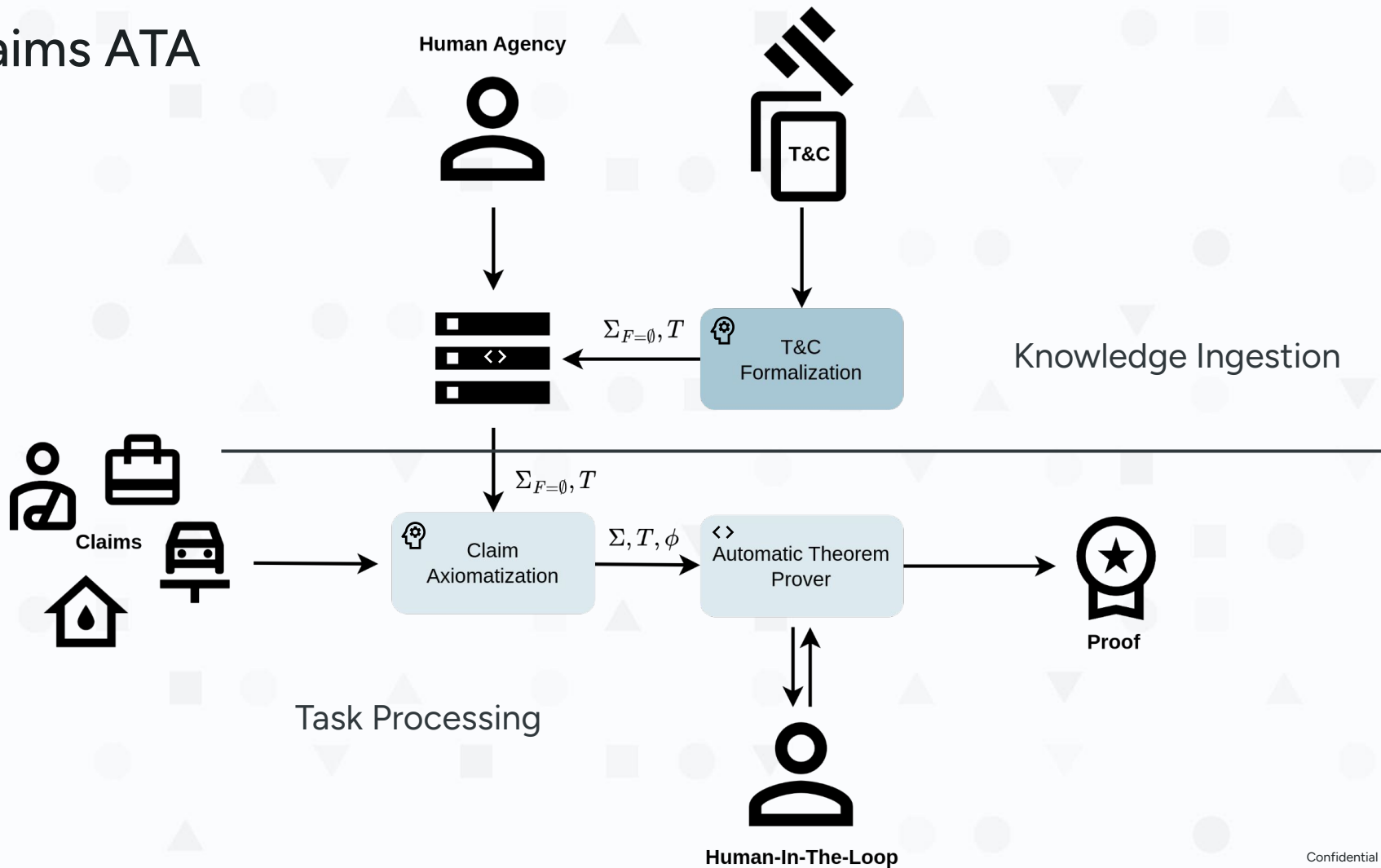
$\neg \exists x R(x,x), R(a,b), \forall xy [R(x,y) \rightarrow x=y \vee R(y,x)] \vdash R(b,a)$

1	$\neg \exists x R(x,x)$	premise
2	$R(a,b)$	premise
3	$\forall xy (R(x,y) \rightarrow x=y \vee R(y,x))$	premise
4	$\forall y (R(a,y) \rightarrow a=y \vee R(y,a))$	$\forall e$ 3
5	$R(a,b) \rightarrow a=b \vee R(b,a)$	$\forall e$ 4
6	$a=b \vee R(b,a)$	$\rightarrow e$ 5, 2
7	$a=b$	assume
8	$R(b,b)$	$=e$ 7, 2
9	$\exists x R(x,x)$	$\exists i$ 8
10	\perp	$\forall e$ 1, 9
11	$R(b,a)$	$\perp e$ 10
12	$R(b,a)$	assume
13	$R(b,a)$	$\forall e$ 7-11 12-13

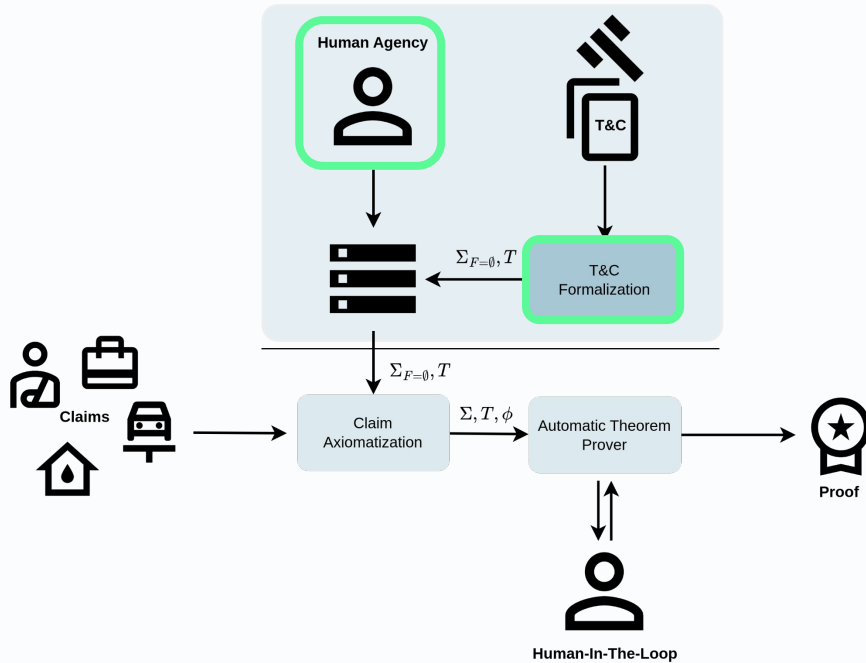
Claims ATA



Claims ATA



Knowledge Ingestion



Predicates & Sorts

`is_relative(p:Person, r:Person):`

Whenever person r is the spouse, partner, mother, sister, ... of person p

...

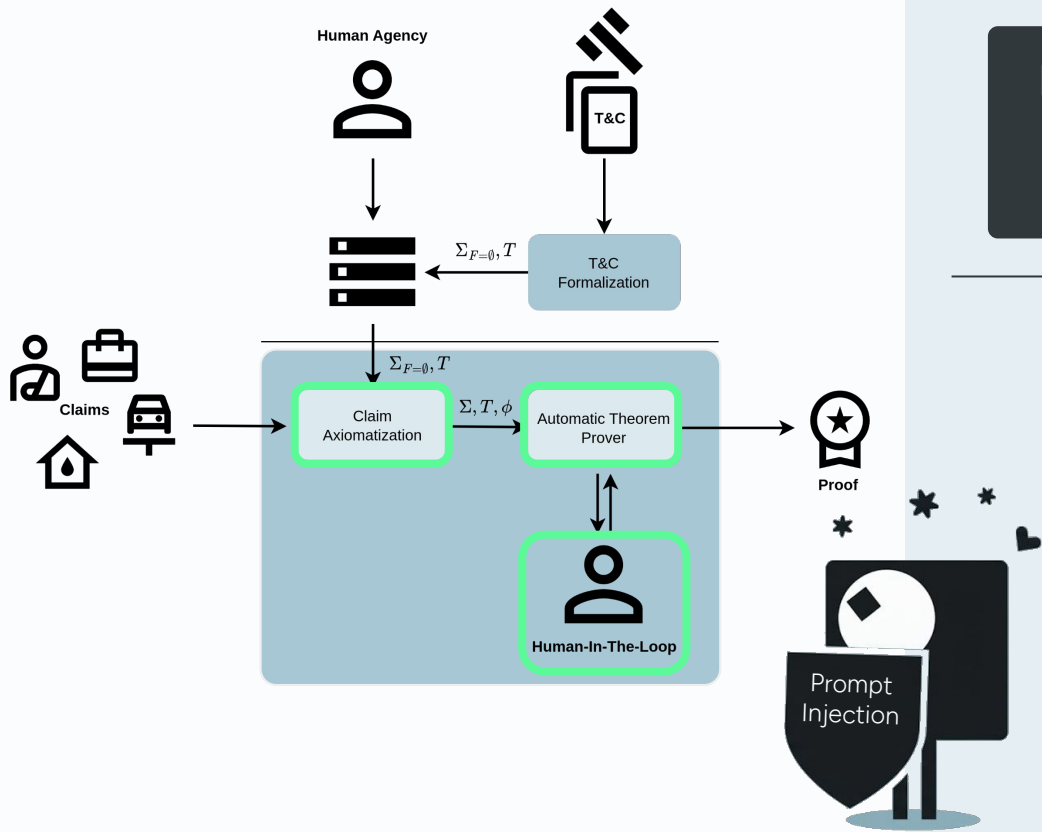
The policyholder is covered for any claim arising from or attributable to sickness or death of (1) you (2) your travelling companion (3) a close relative.

?=

Formula

$$\forall p, r. (\text{is_sick}(p)) \vee (\text{is_dead}(p))$$
$$\vee (\text{is_sick}(r) \wedge \text{is_travelling_with}(r, p))$$
$$\vee (\text{is_dead}(r) \wedge \text{is_travelling_with}(r, p))$$
$$\vee (\text{is_sick}(r) \wedge \text{is_relative}(r, p))$$
$$\vee (\text{is_dead}(r) \wedge \text{is_relative}(r, p))$$
$$\rightarrow \text{is_covered}(p)$$


Task Processing



Alice and her brother Bob booked a trip to Spain. Unfortunately, Bob got the flu...

Named Entity Recognition (NER)
Relation Extraction (RE)
is all you need

`is_covered(ALICE)`

Automatic Theorem Prover

Valid

Invalid



Autonomy, Control, and Velocity for processes where **quality**, **consistency**, and **trust** are non-negotiable.



Accuracy

(Zero Shot)

	Base	ATA
Travel	66%	73%
Electro.	73%	74%
Dental	92%	89%

Stability

(Standard Deviation)

	Base	ATA
Int.	± 1.2	± 0.0
Ext.	± 0.3	± 0.2
LLM	± 4.9	± 1.0

Controllability

(Compare Formula & Rules)

gemini-2.5-pro:	77%
ATA:	73%
<hr/>	
Miss. Clause	+4.7%
Miss. Predicate	+4.8%
Predicate Def.	+4.7%
<hr/>	
ATA + Corr:	87%
with flash-lite:	77%

ATA ... uses gemini-2.5-flash without thinking
Base ... gemini-2.5-flash with built-in thinking



Trustworthy¹ AI

**Accuracy &
Stability**

**Transparency &
Explainability**

**Human Agency &
Oversight**

**Privacy &
Security**

Accountability

**Fairness &
Non-Discrimination**

References



[1] Kowald, Dominik, et al. "*Establishing and evaluating trustworthy AI: overview and research challenges.*" *Frontiers in Big Data* 7 (2024): 1467222.

[2] Peer, David & Stabinger, Sebastian "*ATA: A Neuro-Symbolic Approach to Implement Autonomous and Trustworthy Agents*", arXiv:2510.16381, <https://arxiv.org/abs/2510.16381>



Autonomous Trustworthy Agents (ATA)



Otera

David Peer
david.peer@otera.ai
Machine Learning Researcher

