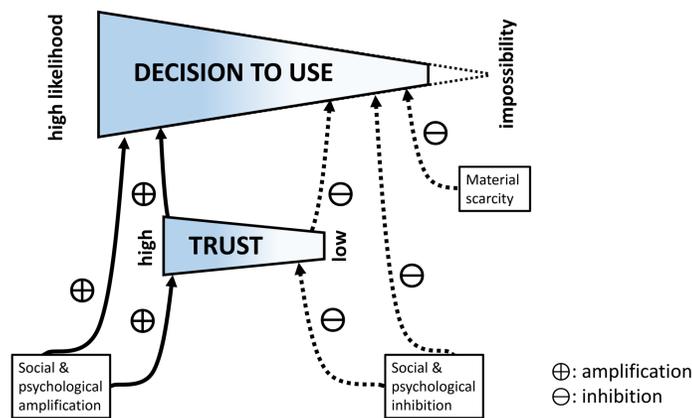


## 1. THE RELEVANCE OF TRUST

Trust is neither a necessary nor a sufficient condition for the decision to rely on a specific technology, but must be seen in relation with other influencing factors and must be contrasted with the notion of *trustworthiness* with which it may be confused. In the following we understand trustworthiness as an objective characteristic reflecting the corresponding levels of *reliability, safety, security, transparency* and *fairness*. Trust, on the other hand, is an attitude based on subjective impression, individual judgment and experience, representing a user's ease of mind in conducting or submitting to an activity. A high level of trust may be considered appropriate when it coincides with a high level of trustworthiness, but trustworthiness does not automatically give rise to trust (resulting in 'undertrust') and trust may be found to be higher than justified when one evaluates the trustworthiness ('overtrust'). Furthermore, trust is by far the only factor affecting the decision to use a device, system or service. Other relevant factors include material availability, psychological dependency, social constraints and personally held beliefs.

likelihood of use	social dependency	psychological dependency	material dependency
very low	ostracization	phobia	material scarcity
	discouragement	aversion	more superior than inferior alternatives
	encouragement	affinity	more inferior than superior alternatives
very high	compulsion	addiction	lack of alternatives



## 2. FRAMEWORK CONCEPTS OF TRUST AND DECISION

- x trusts the means m regarding the achievement of goal g based on the beliefs in
- the capability of m (based on its specifications)
  - the availability of the required functionality (based on its selectable modes of operation)
  - the transparency of m (x believes there are no hidden costs or side effects)
  - the fairness of m (x believes there are e.g. no algorithmic biases involved)
  - goal fulfillment without harm (x believes the likelihood of achieving g through m is high, the risk incurred by not achieving g through m while attempting to do so is low and the risk to others is low (or x chooses to ignore these risks!))
- x then decides to use m IF
- the level of trust T is above a determined level T<sub>min</sub>
- AND
- a subjective appraisal of the total value of goal achievement, costs, risks and side effects leads to a value V greater than a determined level V<sub>min</sub>

## 3. Core Goals and Their Complements

A goal **g** comprises the *core goal* **g\*** (which the user ideally would like to achieve without any costs or side effects) and its *complements* **g\*|<sub>m</sub>** with respect to its achievement, while utilizing the means **m**:

$$g = (g^*, g^*|_m) = (g^*, \text{costs}|_m, \text{side effects}|_m)$$

We call these the *complements of g\* with respect to means m*, because **g\*** cannot be achieved without incurring them other than by changing the means (or possibly the mode) of achievement.

Side effects include risks. With respect to risks we distinguish

- 1<sup>st</sup> person risks (risks to the trustor)
- 2<sup>nd</sup> person risks (risks to persons intentionally participating in the trustor's activity)
- 3<sup>rd</sup> person risks (risks to persons randomly encountered and not intentionally involved)
- n<sup>th</sup> person risks (risks to persons in other locations, who are usually not encountered by and whose existence may even be unknown to the trustor)
- environmental risks

goal **g** = (core goal **g\***, cost **c**, side effects **s**)

subjective allocation of value: VAL

$$VAL(\langle g|m \rangle) = VAL(\langle g^*|m \rangle) - VAL(\langle c|m \rangle) - VAL(\langle s|m \rangle)$$

objective risk determination: RISK

$$RISK(\langle g|m \rangle) = RISK(\langle g^*|m \rangle|x1) \wedge RISK(\langle g^*|m \rangle|x2) \wedge RISK(\langle g^*|m \rangle|x3) \wedge RISK(\langle g^*|m \rangle|xn)$$

$$VAL_{total} = VAL(\langle g|m \rangle) - VAL(RISK(\langle g|m \rangle))$$

## 4. THE UTILITY OF ACTIONS

We can relate an individual action to a) the expectation by the user regarding that action, b) an intrinsic utility of the action and c) the attitude of the user. Expectation contributes to utility through the level of correspondence between action and expectation: the closer the match, the higher the value. If there is a complete mismatch between the two, the utility value assumes a negative sign. The intrinsic utility component is independent of this and expresses the fact that actions which pose contingent (unrelated to goal achievement) objective hazards to the user can never have a positive value. As for the attitude, we consider three states – accepting, rejecting and detached. In the accepting state, the overall resulting utility of an action has the same sign as the result determined by the expectation and intrinsic utility. In the rejecting state, utility is always negative. In the detached state, the utility value is multiplied with the imaginary unit *i*, expressing a situation, where the user's expectations have been frustrated to the degree that only adversarial actions are expected – any "benevolent" action by the autonomous system will then likely be ignored and its (potential) utility, while extant, becoming inconsequential.)

Symbolically, we write:

$$utility_{exp} = \langle \text{Expectation} | \odot | \text{Action} \rangle$$

$$utility_{intrinsic} = \rho; \rho \in \mathbb{R}$$

$$\text{attitude: } \text{sgn } A; \text{sgn } \in \{+, -, i\}; A \in \mathbb{R}_+$$

where ' $\langle \text{Expectation} |$ ' and ' $| \text{Action} \rangle$ ' stand for mathematical representations of expectation and action, to be combined by a suitable operation ' $\odot$ '.

$$|utility_{res}| = |\text{sgn } A| \times |utility_{intrinsic}| \times |utility_{exp}|$$

Whenever any of these three components has a negative sign, the sign of the resulting utility value is negative.

## 5. THE DYNAMICS OF TRUST

Trust at a given step *n* in a sequence of transactions may depend on a) the value *v<sub>n</sub>* of an action at that instance, b) the level of trust at the previous step *T<sub>n-1</sub>*, and c) further characteristics of the overall time series of trust levels  $\chi_{hist}(T_1, \dots, T_{n-1})$ , where  $\chi_{hist}$  is an operator expressing specific properties of the time series, e.g. its variability. These dependencies need to be determined by empirical studies.

The trustor's possible actions (requests with respect to the autonomous system) and state of mind can then be modeled as a labelled state transition system (coupled to a model of the autonomous system), where trust appears as a trace history-dependent label and state transitions depend both on the value of the response by the autonomous system and the level of trust. While the actual value of trust has to be calculated at each step, for the state labelling it is only important whether that value is above or below critical threshold levels.

determine initial value of trust *T<sub>0</sub>*;

determine initial set of expectations *e<sub>0</sub>*;

while use do

```

Tx(ai); /* request action ai */
Rx(ai); /* receive action response */
ui = <ei | ⊙ | ai>; /* determine utility ui */
Ti = f(Ti-1, ui, χhist(T1, ..., Tn-1)); /* calculate trust T */
σ(anext): anext ∈ {a1, ..., an}; /* determine next action request */
enext ← (anext); /* determine next expectation */
if Ti < Tl ∨ Ti > Th then /* if trust out of range */
  if ui < 0 ∨ Ti < Tcritical then
    discontinue use;
  end
  e → e*; /* adjust expectation */
end
end

```

## 6. IMPLICATIONS FOR DESIGN AND VERIFICATION

Systems do fail with a nonvanishing probability, no matter how well designed or verified they are. What one wants to avoid are systems which fail at a high level of user trust.

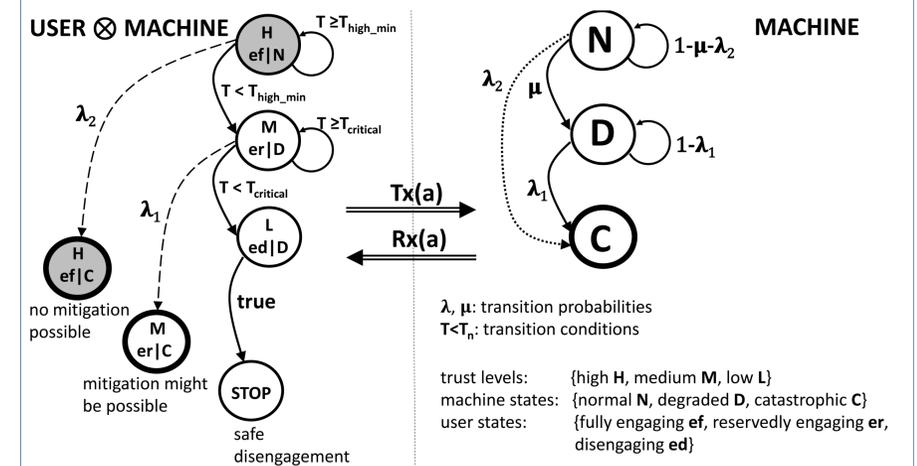
From a verification perspective one wants to assure that

EITHER

the probability of transiting to a critical state at a high level of user trust is low

OR

the level of user trust is sufficiently low for all transitions leading to a critical state.



$\lambda, \mu$ : transition probabilities  
 $T < T_n$ : transition conditions

trust levels: {high H, medium M, low L}  
machine states: {normal N, degraded D, catastrophic C}  
user states: {fully engaging ef, reservedly engaging er, disengaging ed}