

Artificial Intelligence as a System

What Perspective does the Systems Theory Reveal?

first edition, ver. 1.01 (en)

30.12.2022

Igor Furgel
(office@furgel.com)

In the present work we like to apply the systems theory to 'artificial intelligence' (AI) systems. One goal we pursue in doing so is to establish the basic distinguishing criteria between different types of AI systems.

We also like to consider the question, what fundamental prerequisites, from a systems theory perspective, an autonomous technical system shall fulfil so that it would be universally applicable with the same (or stronger and more extensive) intellectual and creative capabilities as a human being.

This consideration will help us to understand the place of different types of Artificial Intelligence systems in this context of 'the inanimate – the animate – the human being'.

The present work may attract the attention of an audience who is interested both in questions of artificial intelligence and its place in the 'the inanimate – the animate – the human being' context, and in the systems approach in general.

The current, first English edition (V. 1.00 (en)) was published on 30.12.2022, Deutsche Nationalbibliothek, <https://d-nb.info/1276876920/>.

Das Original der ersten Ausgabe dieses Aufsatzes (V. 1.00 (de) – „Künstliche Intelligenz als System: Welche Perspektive zeigt die Systemtheorie auf?“) wurde am 30.12.2022 veröffentlicht, Deutsche Nationalbibliothek, <https://d-nb.info/1276877072/>.

Первое издание этой работы на русском (V. 1.00 (ru) – „Искусственный интеллект как система: Какую перспективу открывает теория систем?“) опубликовано 30.12.2022, Deutsche Nationalbibliothek, <https://d-nb.info/1276877293/>.

The main ideas of this work, as applied to AI, emerged between June and December 2022.

Contents

1	Introduction	5
2	Rule-based AI (expert systems)	5
3	Weak AI with Machine Learning	9
3.1	Learning Subsystem of Technical AI System (TSAI-ML)	11
3.2	Learning Subsystem of Automated Decision-Making System (ADMS-ML)	15
3.3	Decision-making Subsystem of Technical AI System (TSAI-DM)	17
3.4	Decision-making Subsystem of Automated Decision-Making System (ADMS-DM)	19
4	Strong AI	22
5	Enmorphya of ‘Self-Awareness’ of AI: Distinction Criteria between Types of Artificial Intelligence and Perspective	26
6	Glossary	29
7	References	35
8	Acknowledgements	35

1 Introduction

In the present work we like to apply the systems theory approach developed in [5], Part A, CHAPTER I (particularly in ch. 3) и CHAPTER II, to ‘artificial intelligence’ (AI) systems. One goal we pursue in doing so is to establish the basic distinguishing criteria between three types of AI systems, namely between

- *rule-based* expert systems, which are based on rule-based programming and were developed in the 1970s and 1980s,
- a ‘Weak AI’¹, which is based on *machine learning* and can only be used to perform specific, limited tasks, and
- a ‘Strong AI’² that would be universally applicable with the same (or stronger and more extensive) intellectual and creative capabilities as a human being; the ‘Strong AI’ is currently only a vision.

We also like to consider the question, what fundamental prerequisites, from a systems theory perspective, an autonomous technical system shall fulfil so that it would be universally applicable with the same (or stronger and more extensive) intellectual and creative capabilities as a human being.

This consideration will help us to understand the place of different types of Artificial Intelligence systems in this context of ‘the inanimate – the animate – the human being’.

This work can be read on its own. The fundamental concepts are given in ch. 6 “Glossary“. Nevertheless, since the approaches we have developed in [5], Part A, CHAPTER I (in particular in ch. 3 “Being, Existential Triads and Enmorphya”) and CHAPTER II are fundamental to this study, we recommend that readers interested in the underlying developments also refer to [5].

2 Rule-based AI (expert systems)

Rule-based expert systems are created using traditional programming, also known as rule-based programming. Such expert systems use hard-coded knowledge bases on a particular subject area and hard-coded decision rules. When a rule-based expert system receives a query with data to be analysed, it uses the hard-coded decision rules to compare this data to be analysed with the knowledge base and generates a response (decision/recommendation), see [4], ch.1 and [3], ch. 2.

Let us consider how the abstract elements of the *existential triad* (see Glossary) are expressed within this type of AI.

Data/queries to be analysed are the ‘*substrate*’ of any *rule-based* expert system. The ‘*property*’ of the system is the property of the rule-based programme (including the operationalisation/measurability of the purpose of the programme, the model of the problem and the algorithm)³, which is hard-coded in programming instructions including the

¹ Narrow (weak) artificial intellect, abbreviated by ANI; DE: Schwache künstliche Intelligenz (Schwache KI); RU: слабый (узкий) искусственный интеллект (ИИ)

² General AI, abbreviated by AGI; DE: Starke KI; RU: сильный (общий) ИИ

³ These terms are explained in detail in [2], see e.g. ‘Glossary’ there. For the sake of better readability, we only give the respective core of the definitions here:

knowledge base, and the properties/characteristics of the data to be analysed. The ‘*relation*’ is the process of applying the hard-coded instructions to the data to be analysed by the programme execution in the context of using the rule-based expert system. The data to be analysed does not change as a result of this application of the instructions.

The process of applying the hard-coded instructions to the data to be analysed is *deterministic*, not *stochastic*. Therefore, ‘the Principle of Sufficiency of the Existential Triad’⁴ is not applicable here: the rule-based expert systems are *deterministic*.

For *rule-based* expert systems, implemented *principles of software development* including the applied procedural rules/norms (i.e. the ‘*programming manual*’) represent the ‘relation-control-information’ (the *enmorphya of relation*). As the enmorphya of the relation between the substrate (the data to be analysed) and the property (the nature of the rule-based programme), the implemented principles of software development (i.e. the ‘*programming manual*’) define the character of this relation (interaction).

If we now draw a parallel with *stochastic* systems, for which ‘the Principle of Sufficiency of the Existential Triad’⁴ generally applies, we find that this principle is replaced by the ‘*programming manual*’ when it deals with a *deterministic* expert system.

The ‘*programming manual*’ determines the character of the application of the hard-coded instructions to the data to be analysed. The ‘*programming manual*’ thereby shapes various *optimisation functions* such as minimising the power consumption or the execution time of the programme. Unlike *stochastic* systems in general, the ‘*programming manual*’ *cannot* shape the data to be analysed, but only the properties of the rule-based programme (including the operationalisation/measurability of the purpose of the programme, the model of the problem and the algorithm).

Let us illustrate the relation between the primary system and the metasystem using the example of rule-based expert systems:

- “By modelling I mean any form of simplification and abstraction of a situation that is still accurate enough to allow predictions or analytical conclusions about that situation.” (In the original: „Unter einer Modellierung verstehe ich jegliche Form der Vereinfachung und Abstraktion einer Situation, die trotzdem noch so genau ist, dass sie Vorhersagen oder analytische Schlüsse über diese Situation zulässt.“)

- “An operationalisation represents the measurability of a (social) concept – it is always based on a model of the concept.” (In the original: „Eine Operationalisierung stellt die Messbarmachung eines (sozialen) Konzeptes dar – sie basiert immer auf einem Modell des Konzeptes.“)

- “An algorithm is a sufficiently detailed and systematic set of instructions for any experienced programmer to solve a mathematical problem so that, if implemented correctly (translated into code), the computer will calculate the correct output for any correct set of inputs.” (In the original: „Ein Algorithmus ist eine für jede erfahrene Programmiererin und jeden erfahrenen Programmierer ausreichend detaillierte und systematische Handlungsanweisung, um ein mathematisches Problem zu lösen, sodass bei korrekter Implementierung (Übersetzung in Code) der Computer für jede korrekte Inputmenge den korrekten Output berechnet.“)

⁴ [5], CHAPTER I, STM. 9 ‘The principle of sufficiency of the existential triad’: “If ‘relation’ in an existential triad {substrate, property, relation} has a fundamentally *stochastic* character and *statistically* obeys a certain law, then this existential triad is not only necessary, but also sufficient for the achievement of observability and, thus, for creating the state of ‘being’ of the system based on this existential triad. The evolution of this system will follow the character of the ‘relation’ in the existential triad.”

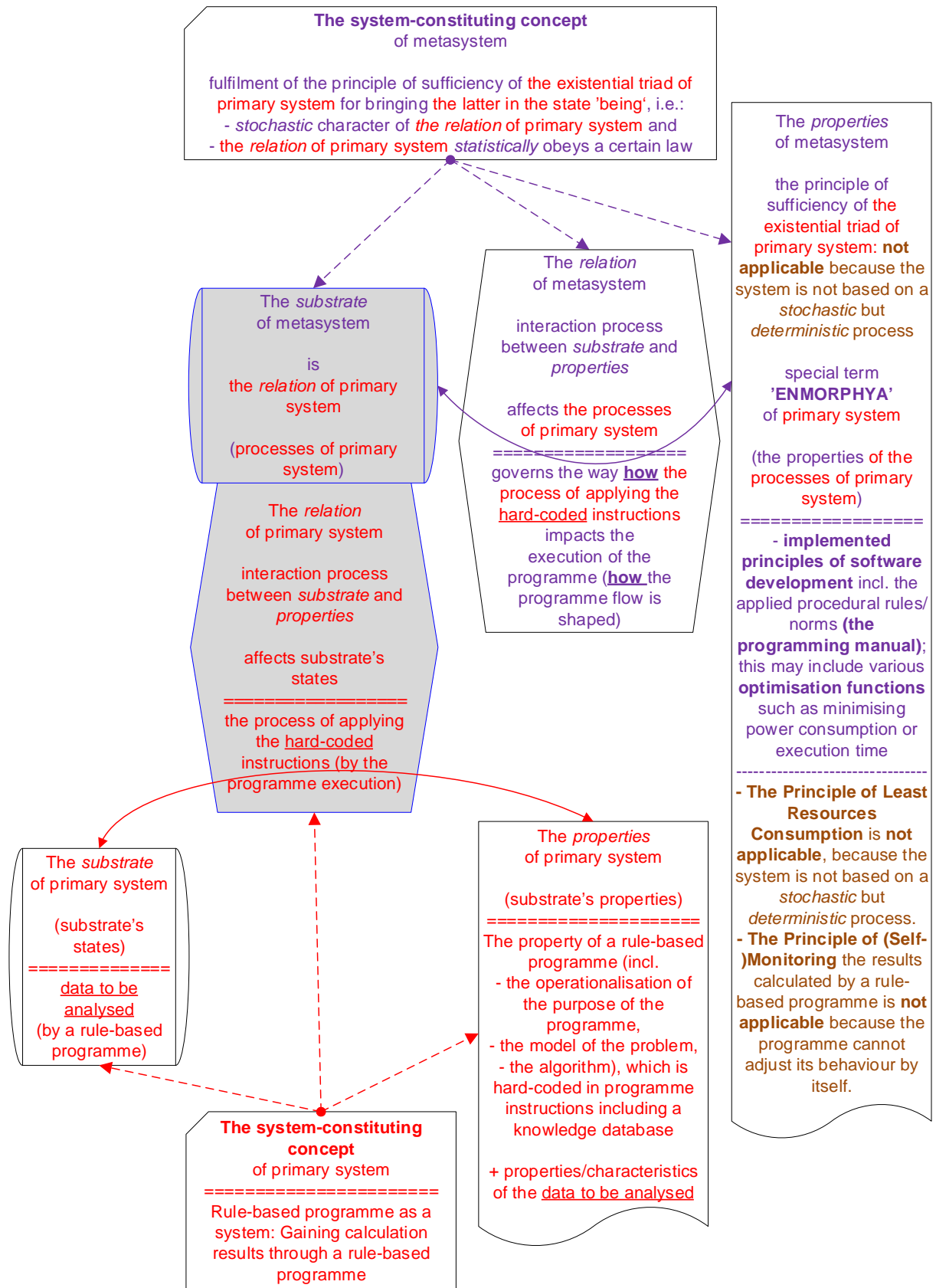


Figure 1: Relationship between the primary system 'rule-based expert system' and the corresponding metasytem

The systems-theoretical consideration of *rule-based* expert systems allows us to make a useful generalisation of the systems-theoretical difference between *deterministic* and *stochastic* systems:

System type →	<i>stochastic</i> systems	<i>deterministic</i> systems (can exist only as <u>artefacts</u> ⁵)
Parameter ↓		
the Principle of Sufficiency of the Existential Triad ⁴	always applies	not applicable
Enmorphya of relation ('relation-control-information')	is a system-specific implementation of 'the principle of sufficiency of the existential triad', whereby the <i>Principle of Least Resources Consumption</i> is always a component of enmorphya ⁶ .	is the practically implemented <i>principles of the manufacturing</i> of the respective system incl. the applied procedural rules/standards; is the <i>technical specification</i> (or, pictorially speaking, the ' <i>programming manual</i> ') for the manufacturing of the respective system.
Impact of the enmorphya of the relation on the primary system	The respective system-specific implementation of 'the principle of sufficiency of the existential triad' determines the character of the 'relation' of the primary system. In this way, the respective system-specific implementation of 'the principle of sufficiency of the existential triad' shapes <u>both</u> the states of the 'substrate' and the 'properties' of the primary system.	The practically implemented <i>principles of the manufacturing</i> of the respective system <u>cannot</u> shape states of the 'substrate', but only the 'properties' of the primary system.
Impact of the 'relation' of	The interaction process of	The interaction process of

⁵ see [5], CHAPTER I, ch. 3.1 „Being and Existential Triads“ or CHAPTER VII, ch. 2.1.3 „Indeterminacy and Action Quanta: Complementary Characters of the Past and the Future“

⁶ see [5], CHAPTER I, ch. 3.2 „Enmorphya“, STM. 11 'the Principle of Least Resources Consumption': "the **Principle of Least Resources Consumption (PLR)** is relation-control-information (i.e. enmorphya of relation) and governs not only the process of interaction between matter and information in nature, but also between the *substrate* and the *structural factor* of any system – physical, social, communicative, etc. – based on a *stochastic process*."

System type →	<i>stochastic</i> systems	<i>deterministic</i> systems (can exist only as <u>artefacts</u> ⁵)
Parameter ↓		
the primary system on states of the ‘substrate’	the ‘properties’ of the primary system with the ‘substrate’ of the primary system <u>changes</u> states of the ‘substrate’.	the ‘properties’ of the primary system with the ‘substrate’ of the primary system <u>cannot</u> change states of the ‘substrate’.

The **Principle of Least Resources Consumption**, which applies to all *stochastic* systems, is not applicable to a *rule-based* expert system because the system is based not on a *stochastic* but on a *deterministic* process⁶.

The **Principle of Self-Preservation of the System**, which applies to all *quasi-stochastic* systems⁷, and which for many such systems is expressed as the Principle of (self)monitoring, is also not applicable because the system is based not on a *stochastic*, but on a *deterministic* process, and the programme cannot therefore adapt its behaviour itself. Based on a result calculated by a *rule-based* expert system, the expert system cannot autonomously adjust its behaviour.

3 Weak AI with Machine Learning

‘Weak AI’ – in contrast to rule-based expert systems – is based on *machine learning* (ML). ‘Weak AI’ can only be used for specific, limited tasks, see [2], Glossary, [3], ch. 1, [4], ch. 2 (“ANI”) or on the internet.

The *machine learning* approach differs technologically from the *rule-based* programming approach in that decision rules and knowledge bases are no longer hard-coded into the programme, as in the case of *rule-based* AI, but the programme implementing the *machine learning* determines a set of decision rules itself in the form of a *statistical model*, based on data fed to the programme with ML⁸.

There are several technological types of machine learning, such as ‘deep learning’ based on artificial neural network (ANN) technology, ‘random forest’ based on decision tree architecture, ‘support vector machine (SVM)’⁹ and some others, see [4], ch. 3, [2], ch. 5.

There are also some organisational types of machine learning, whereby organisational and technological types of machine learning should fit together so that ‘Weak AI’ can efficiently achieve the predefined purpose. As examples of the organisational types of machine learning, we mention here (see [3], ch. 4):

⁷ see [5], CHAPTER I, ch. 3.5 „Enmophya for Quasi-Stochastic Systems“, STM. 13 ‘the Principle of Self-Preservation of System’: “in order to ensure the stability of *quasi-stochastic* systems, their enmophya shall contain at least one more principle, which we called the **Principle of Self-Preservation of System**.”

⁸ Machine learning (ML) is defined in [2], Glossary as follows: “A collection of methods that look for patterns in data from the past that allow predictions for the future.” (In the original: “Eine Sammlung von Methoden, die in Daten der Vergangenheit nach Mustern suchen, die für die Zukunft Vorhersagen erlauben.”)

⁹ RU: машина опорных векторов

- Supervised learning¹⁰: The programme with ML is first fed *training data* selected by the human (in the role of the ‘data scientist’, see [2]), on the basis of which the programme with ML searches for correlations (and usually finds them) and on this basis determines a set of decision rules in the form of a statistical model. The programme with ML is then fed *test data* (i.e. data with already known results¹¹). The ‘Weak AI’ analyses this test data using the decision rules it set itself during training and returns a result of this analysis. The human (in the role of the ‘data scientist’) compares the previously known result of the test data with the returned result of the ‘weak AI’ and applies a *quality measure*¹² defined by the human to this comparison. If the deviation between the compared results satisfies this *quality measure*, the training is considered successful, cf. [2], ch. 5, [4], ch. 3 (“Neuronale Netzwerke“ (‘Neural Networks’)).
- - Unsupervised learning¹³: The programme with ML is fed unlabelled data, i.e. without any predefinition of the learning goal. The programme with ML searches in this data set for correlations, i.e. for dependencies and patterns in these data and derives a set of decision rules from the correlations¹⁴ found. If this programme with ML is fed further data for analysis, the programme with ML continues the search for further correlations, readjusts the decision rules and applies them to the input data. As a result of this application, the input data are assigned to the clusters (categories) learned in this way. One of the variants of unsupervised learning is ‘Generative Adversarial Networks (GAN)’. Since there is no ex-ante ‘ground truth’¹¹ in unsupervised learning, the human must check ex-post whether the decisions/recommendations generated as a result by this type of ‘Weak AI’ are adequate for the objective of the operator of the AI system, see [3], ch. 4.
- Reinforcement learning¹⁵: In the first step, the human (in the role of ‘data scientist’) defines a *reward function* for the programme with ML and instructs the programme to maximise the value of the ‘*reward*’ when searching for the possible results of the analysis of the input data. The ‘Weak AI’ then tries possible solution options (trial and error) and picks the option that increases the value of the ‘*reward*’ compared to the previous solution as the appropriate decision. In this way, the ‘Weak AI’ optimises itself further and further with each next solution, so that the value of the ‘*reward*’ increases on a statistical average and at some point reaches the maximum value, see [3], ch. 4, [4], ch. 3 (“Evolutionary Algorithms” (“Evolutionäre Algorithmen”)).

One can see from the examples above that a ‘Weak AI’ is dependent on ex-ante or ex-post human intervention. This fact also applies to all other organisational types of machine learning.

Therefore, systems based on machine learning are always *socio-technical* systems. This means that they necessarily have a purely technical core, i.e. a purely technical subsystem, which is embedded in a well-organised *human environment*, see [2], ch. 1, figure 3.

¹⁰ DE: Überwachtes Lernen; RU: обучение под наблюдением (с учителем; контролируемое обучение)

¹¹ The characteristics of the *training* and *test data*, which are already known in advance, are called ‘ground truth’.

¹² “Quality measure: A function that evaluates how good an (algorithmic) solution of a problem is.” (In the original: „Qualitätsmaß: Eine Funktion, die bewertet, wie gut eine (algorithmische) Lösung eines Problems ist.“), see [2], Glossary.

¹³ DE: Unüberwachtes Lernen; RU: обучение без наблюдения (без учителя; неконтролируемое обучение)

¹⁴ It is important to note that correlations found are not necessarily causal, what is something a ‘Weak AI’ fundamentally cannot recognise.

¹⁵ DE: Bestärkendes Lernen; RU: обучение с подкреплением

The corresponding entire socio-technical system – i.e. the purely technical subsystem together with the organised human environment – is called an ‘*Automated (Algorithmic) Decision-Making System*’ (ADMS¹⁶), see [2], ch. 1, [3], ch. 1.

The purely technical subsystem (we will abbreviate it as TSAI – *technical system of artificial intelligence*) and the entire ADMS (*automated decision-making system*) belong – from a systems-theoretical point of view – to different systems categories (more on this – later in this section). Therefore, it is necessary to systems-theoretically distinguish between TSAI and ADMS, see [2], ch. 1, figure 2.

Another difference relevant from a systems-theoretical point of view exists within each TSAI. The architecture of each TSAI comprises two technical subsystems: The machine learning subsystem (the ML subsystem), which functions *stochastically*, and the decision-making subsystem, which functions *deterministically* based on the set of decision rules defined by the ML subsystem, see [2], ch. 5, figure 22.

Accordingly, we have to distinguish – considered from a systems theory perspective – between the following subsystems of a ‘Weak AI’:

- TSAI-ML: Learning subsystem of technical AI,
- ADMS-ML: Learning subsystem of automated decision-making system,
- TSAI-DM: Decision-making subsystem of technical AI, and
- ADMS-DM: Decision-making subsystem of automated decision-making system.

Below we look at each of these subsystems separately.

3.1 Learning Subsystem of Technical AI System (TSAI-ML)

Decision rules to be learned (the statistical model) of the TSAI are the ‘substrate’ of a TSAI-ML subsystem of a ‘Weak AI’. The ‘property’ of the subsystem is the property of the TSAI (including the operationalisation/measurability of the purpose of the TSAI, the model of the problem and the algorithm)³ and the properties/characteristics of the *training data*. The ‘relation’ is the interaction process of the properties of the TSAI and the characteristics of the training data on the one hand with the decision rules to be learned on the other hand, i.e. the actual learning process.

The operationalisation of the purpose of the TSAI includes, among other things, the *fairness measure*¹⁷, which, if relevant, is defined by humans (in the role of the data scientist), e.g. equality vs. equity, see [2], ch. 8.

Let us consider the ‘deep learning’ technology as an example. In this case, the ANN (artificial neural network) starts the learning process by setting the weights/probabilities of the transitions between the single states of the neighbouring layers of artificial ‘neurons’ to random values by the ANN. After each learning cycle using *training data* (as part of the *ground truth*¹¹), the ANN *probabilistically* readjusts these weights with the aim of bringing its

¹⁶ DE: Automatisiertes (Algorithmisches) Entscheidungssystem; RU: автоматизированная система принятия решений

¹⁷ “Fairness measure: A mathematical function that assesses the extent to which different groups of the population are equally or equitably affected by decisions.” (In the original: „Fairnessmaß: Eine mathematische Funktion, die bewertet, inwieweit unterschiedliche Bevölkerungsgruppen gleichermaßen von Entscheidungen betroffen sind“), see [2], Glossary.

decision result closer to the ground truth¹⁸. Accordingly, the sequence of sets of decision rules to be learned that make up the statistical model is also fundamentally *probabilistic*.

Since each next state of the set of decision rules to be learned depends probabilistically solely on its current state (and not on previous states), the learning process has the *Markov property*. Therefore, the learning process is a *true-stochastic* process, see Glossary.

As we stated in [5], CHAPTER I, ch. 3.4 „Enmorphya for Truly-Stochastic Systems“, the ‘relation-control-information’ (enmorphya of relation) of all *truly-stochastic* systems is always represented by the Principle of Least Resources Consumption (of maximum entropy). This means that the enmorphya of the relation of the TSAI-ML subsystem must also be represented by the Principle of Least Resources Consumption (PLR), specifically – by the *Principle of learning economy*.

Manuela Lenzen notes in [3], ch. 4:

“A good ANN is confident in assigning data to the desired categories. It is neither too sensitive nor too insensitive to variations in the data. And it is thrifty in terms of the time, the amount of data and the hardware needed to train it.”¹⁹

Janelle Shane quotes Alex Irpan, AI researcher at Google, in [4], ch. 5:

“I’ve taken to imagining deep RL as a demon that’s deliberately misinterpreting your reward and actively searching for the laziest possible local optima. It’s a bit ridiculous, but I’ve found it’s actually a productive mindset to have.”²⁰

In view of the insight that the TSAI-ML subsystem of a ‘Weak AI’ always follows the *Principle of learning economy* (of the least resources consumption), this statement by Alex Irpan not only no longer sounds “ridiculous”, but downright consequential.

In [4], ch. 6, Janelle Shane repeatedly notes that during training an AI keeps trying to hack the given ‘matrix’ in order to get ‘free energy/food’.

This is also immediately explainable from the *Principle of learning economy*: Hacking the ‘matrix’ always minimises resources consumption.

The second principle of the enmorphya of the relation of the TSAI-ML subsystem is a set of *optimisation functions* according to the purpose of the TSAI including their prioritisation, whereby the *quality measure* and *fairness measure* are specified outside the TSAI, i.e. externally to the TSAI, by the data scientist.

Within the TSAI-ML subsystem, the *Principle of learning economy* together with *optimisation functions* constitute the ‘relation-control-information’ (enmorphya of relation) of this system. As the enmorphya of the relation between the substrate (the decision rules to be learned) and the property (the property of the TSAI and the characteristics of the *training data*), the *Principle of learning economy* and the *optimisation functions* define the character of this relation (interaction), see footnote 4. They determine the character of the machine

¹⁸ The quality of the decision result is determined by the human-defined *quality measure*.

¹⁹ „Ein gutes KNN ist sicher darin, Daten den gewünschten Kategorien zuzuordnen. Es ist weder zu empfindlich noch zu unempfindlich gegenüber Variationen in den Daten. Und es ist sparsam, was die Zeit, die Menge an Daten und die Hardware angeht, die man zum Training benötigt.“

²⁰ Alex Irpan *Deep Reinforcement Learning Doesn’t Work Yet*, <https://www.alexirpan.com/2018/02/14/rl-hard.html> (In [4]: “Ich habe mir angewöhnt, [die KI] als einen Dämon zu betrachten, der seine Belohnung absichtlich falsch interpretiert und aktiv nach dem Optimum sucht, bei dem sie möglichst faul sein kann. Das klingt lächerlich, aber die Einstellung kann tatsächlich recht produktiv sein.”)

learning process, which in turn implements the interaction between the decision rules to be learned and the property of the TSAI and the characteristics of the *training data*. The *Principle of learning economy* and the *optimisation functions* thus shape both the decision rules to be learned (the substrate of the TSAI-ML subsystem) and the property of the TSAI and the matching characteristics of the *training data* (their form and content).

Let us illustrate the relation between the primary system and the metasystem using the example of TSAI-ML subsystem (learning subsystem of technical AI):

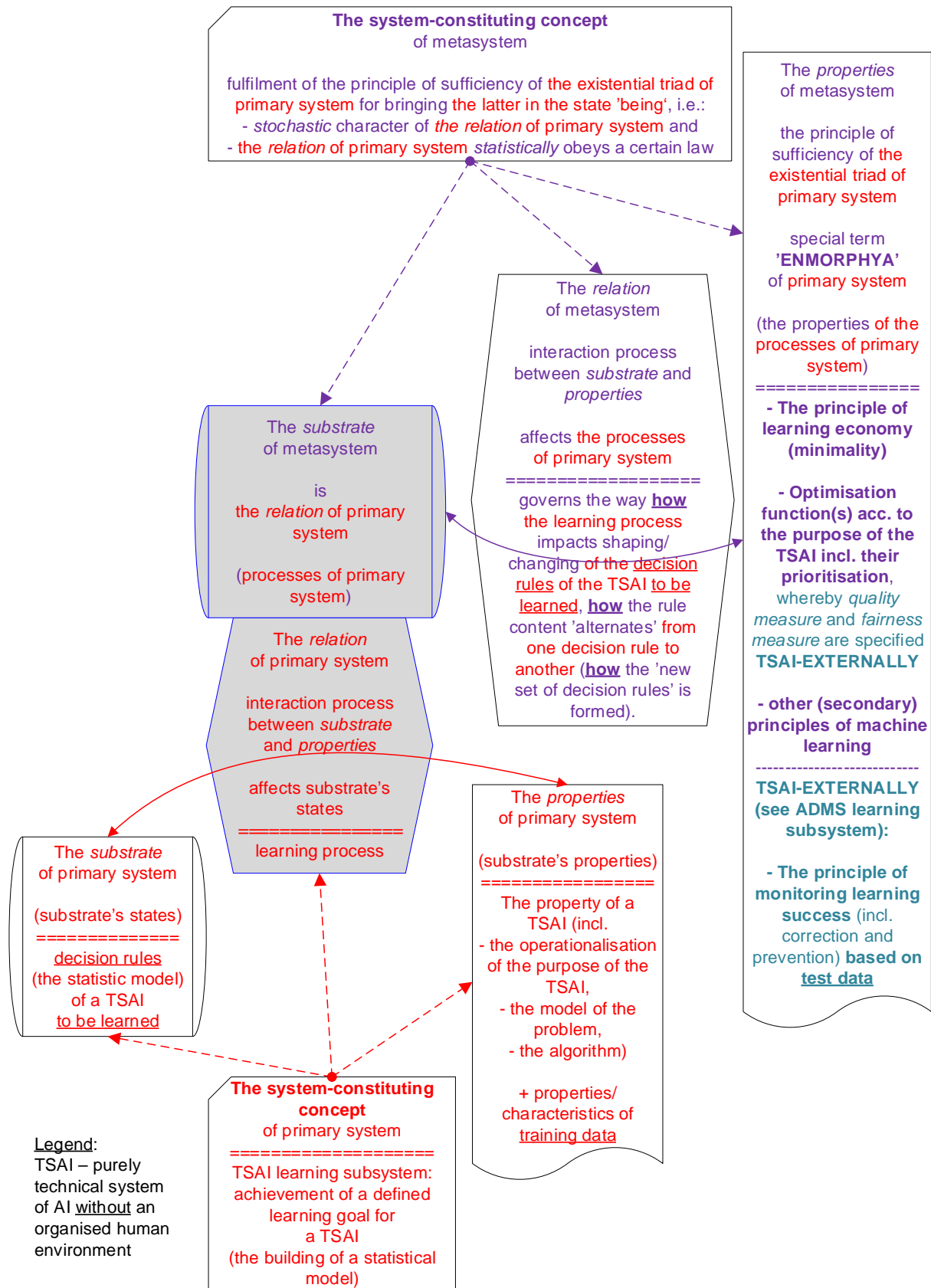


Figure 2: Relationship This means, among other things, that a 'Weak AI' as a purely between the primary system TSAI-ML subsystem (technical AI system – the learning subsystem) and the corresponding metasytem

3.2 Learning Subsystem of Automated Decision-Making System (ADMS-ML)

As we have already described in this chapter, ‘Weak AI’ is always dependent on ex-ante or ex-post human intervention. That is why systems based on machine learning are always *socio-technical* systems. This means that they necessarily have a purely technical core, i.e. a purely technical subsystem (the TSAI), which is embedded in a well-organised human environment, see [2], Ch. 1, Fig. 3. The corresponding overall socio-technical system – i.e. the purely technical subsystem (the TSAI) together with the organised human environment – is called ‘Automated (Algorithmic) Decision-Making System’ (ADMS), see [2], ch. 1, [3], ch. 1.

In ch. 3.1 above, we considered the learning subsystem of technical AI (the TSAI-ML subsystem) from a systems theory perspective. We now want to analyse the learning subsystem of the entire automated decision system (the ADMS-ML subsystem).

From the point of view of the purely technical TSAI, the ‘socio’ component, i.e. human intervention, represents a boundary condition external to the system. The learning subsystem of the entire automated decision-making system differs from the learning subsystem of the technical AI only in that it also integrates the ‘socio’ component, i.e. human intervention.

On Figure 2 it is easy to see that human intervention – as a boundary condition external to the system for the TSAI-ML subsystem – exclusively influences the enmorphya of the relation of the TSAI-ML subsystem (marked by turquoise). The human (in the role of the data scientist) specifies the *quality measure* and the *fairness measure* (as well as other necessary *hyperparameters*) and monitors the learning success on the basis of *test data*, see [2], ch. 5.

The monitoring of learning success (incl. *correction* and *prevention*) based on *test data* implements the *principle of monitoring of learning success* for the learning subsystem of the Automated Decision-Making System (ADMS-ML subsystem), cf. [2], ch. 5, [4], ch. 3. The *principle of monitoring* incl. *correction* and *prevention* is implemented by the system’s own *adaptation mechanism*.

The *principle of monitoring of learning success* is in turn the concrete implementation for the ADMS-ML subsystem of the general **Principle of Self-Preservation of System**, see footnote 7. The fact that the enmorphya of the relation of the ADMS-ML subsystem includes the Principle of Self-Preservation of System indicates that the ADMS-ML subsystem is a *quasi-stochastic* system, see footnotes 7 and 21. This is perfectly comprehensible, however, because the human being, which is itself a *quasi-stochastic* system²², is a component of the ADMS-ML subsystem.

Let us illustrate the relation between the primary system and the metasystem using the example of ADMS-ML subsystem (learning subsystem of automated decision-making system):

²¹ see [5], CHAPTER I, ch. 3.5 „Enmorphya for Quasi-Stochastic Systems“, STM. 14 ‘at least two principles are existentially necessary components of the enmorphya’: “at least two principles are existentially necessary components of the enmorphya of *quasi-stochastic* systems: **the Principle of Least Resources Consumption (PLR)** and **the Principle of Self-Preservation of System (PSP)**.”

²² see [5], CHAPTER II, ch. 1.2 “Principles of the Enmorphya of Living Beings”

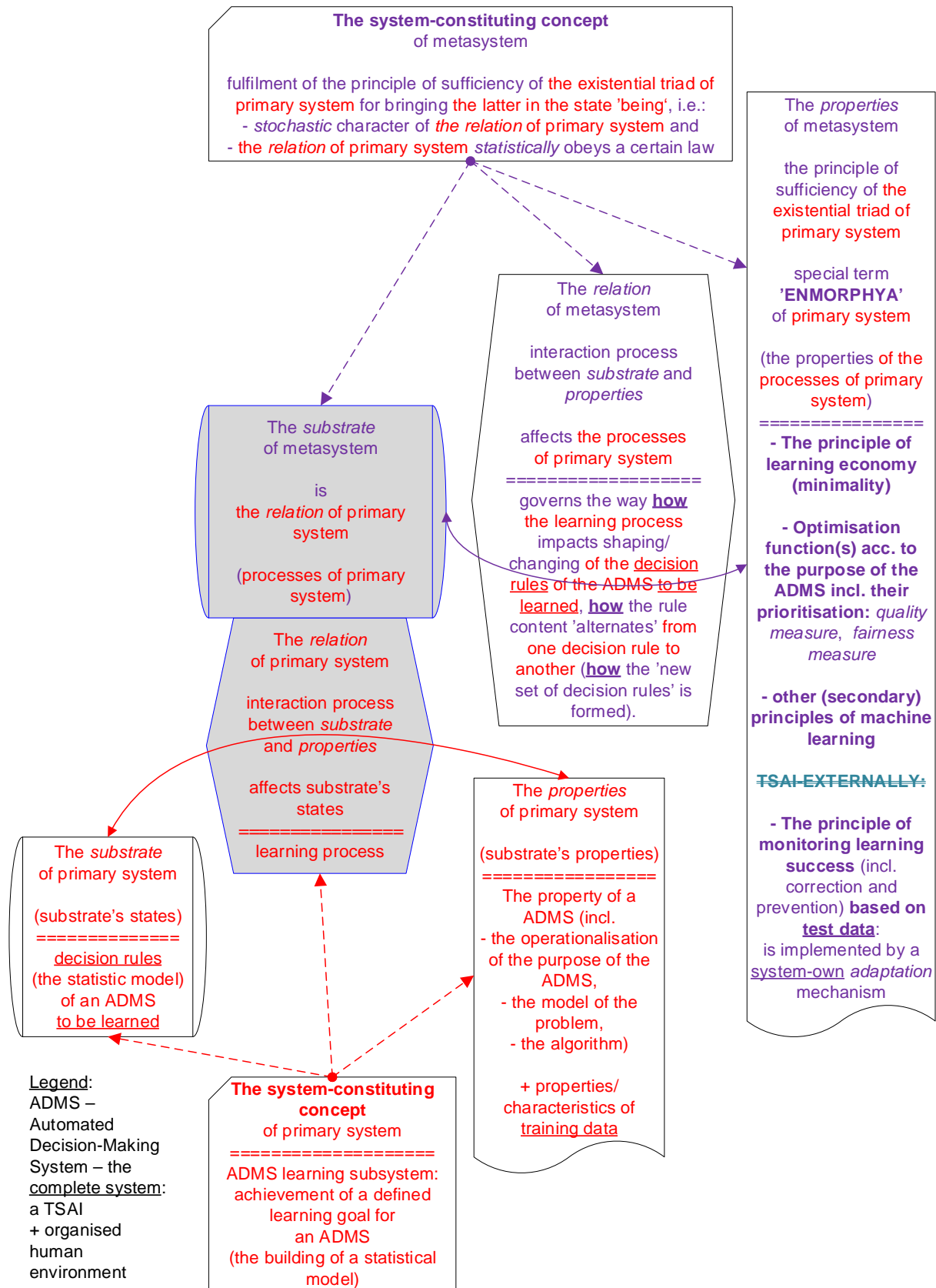


Figure 3: Relationship between the primary system ‘ADMS-ML subsystem’ (automated decision-making system – the learning subsystem) and the corresponding metasytem

3.3 Decision-making Subsystem of Technical AI System (TSAI-DM)

Data to be analysed is the ‘substrate’ for each TSAI decision-making subsystem of a ‘Weak AI’. The ‘property’ of the subsystem is the learned decision rules (the learned statistical model) of the TSAI and the properties/characteristics of the data to be analysed. The ‘relation’ is the process of applying the learned decision rules to the data to be analysed. The data to be analysed does not change through this application.

The process of applying the learned decision rules to the data to be analysed is *deterministic*, not *stochastic*. Therefore, ‘the Principle of Sufficiency of the Existential Triad’⁴ is not applicable here: the TSAI decision-making subsystem is *deterministic*, cf. [2], ch. 5, figure 22.

For the TSAI decision-making subsystem, the *principles of the learned statistical model* represent the relation-control-information (the enmorphya of the relation). The *principles of the learned statistical model* are defined by the TSAI learning subsystem, see “enmorphya” on Figure 2. As the enmorphya of the relation between the substrate (the data to be analysed) and the property (the learned statistical model), the principles of the learned statistical model define the character of this relation (interaction). Comparing the enmorphya of the rule-based expert system (Figure 1) and the enmorphya of the TSAI decision-making subsystem (Figure 4 below), it becomes obvious that the learned statistical model for a ‘Weak AI’ plays the same (central) role as the ‘programming manual’ for rule-based expert systems.

Let us illustrate the relation between the primary system and the metasystem using the example of TSAI-DM subsystem (decision-making subsystem of technical AI system learning):

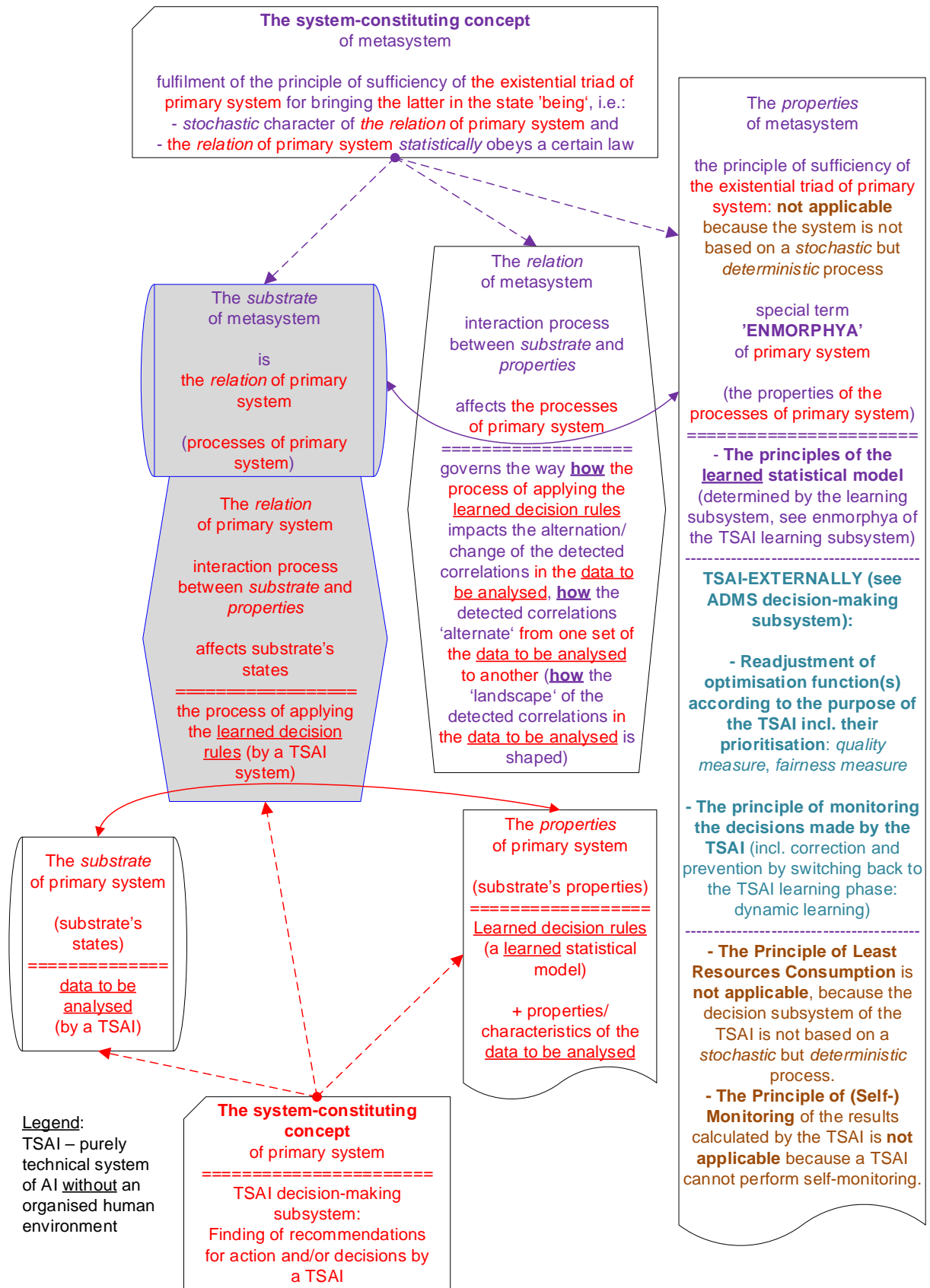


Figure 4: Relationship between the primary system TSAI-DM subsystem (technical AI system – the decision-making subsystem) and the corresponding metasytem

3.4 Decision-making Subsystem of Automated Decision-Making System (ADMS-DM)

In ch. 3.3 above, we looked at the decision-making subsystem of technical AI (the TSAI decision-making subsystem) from a systems theory perspective. We now want to analyse the decision-making subsystem of the entire *socio-technical* automated decision-making system (the ADMS decision-making subsystem).

From the point of view of the purely technical TSAI, the ‘socio’ component, i.e. human intervention, represents a boundary condition external to the system. The learning subsystem of the entire automated decision-making system differs from the learning subsystem of the technical AI only in that it also integrates the ‘socio’ component, i.e. human intervention.

On Figure 4 it is easy to see that human intervention – as a boundary condition external to the system for the TSAI-DM subsystem – exclusively influences the enmorphya of the relation of the TSAI-DM subsystem (marked by turquoise). The human (in the role of the data scientist) readjusts *optimisation function(s)* according to the purpose of the ADMS including their prioritisation as well as the *quality measure* and the *fairness measure* as well as other necessary *hyperparameters*. Besides, the human monitors the adequacy of decisions made, cf. [2], ch. 8 and 9, [4], ch. 3. This implements the *principle of monitoring the decisions made by the ADM system* (incl. *correction* and *prevention* by switching back to the ADMS learning phase, Figure 3: dynamic learning).

The *principle of monitoring* incl. *correction* and *prevention* is implemented by the system’s own *adaptation mechanism*.

Similar to the ADMS-ML subsystem (see ch. 3.2), the *principle of monitoring the decisions made* is in turn the concrete implementation for the ADMS-DM subsystem of the general **Principle of Self-Preservation of System**, see footnote 7. The fact that the enmorphya of the relation of the ADMS-DM subsystem includes the Principle of Self-Preservation of System indicates that the ADMS-DM subsystem is a *quasi-stochastic* system, see footnotes 7 and 21. This is perfectly comprehensible, however, because the human being, which is itself a *quasi-stochastic* system²³, is a component of the ADMS-DM subsystem.

Let us illustrate the relation between the primary system and the metasystem using the example of ADMS-DM subsystem (decision-making subsystem of automated decision-making system):

²³ see [5], CHAPTER II, ch. 1.2 “Principles of the Enmorphya of Living Beings”

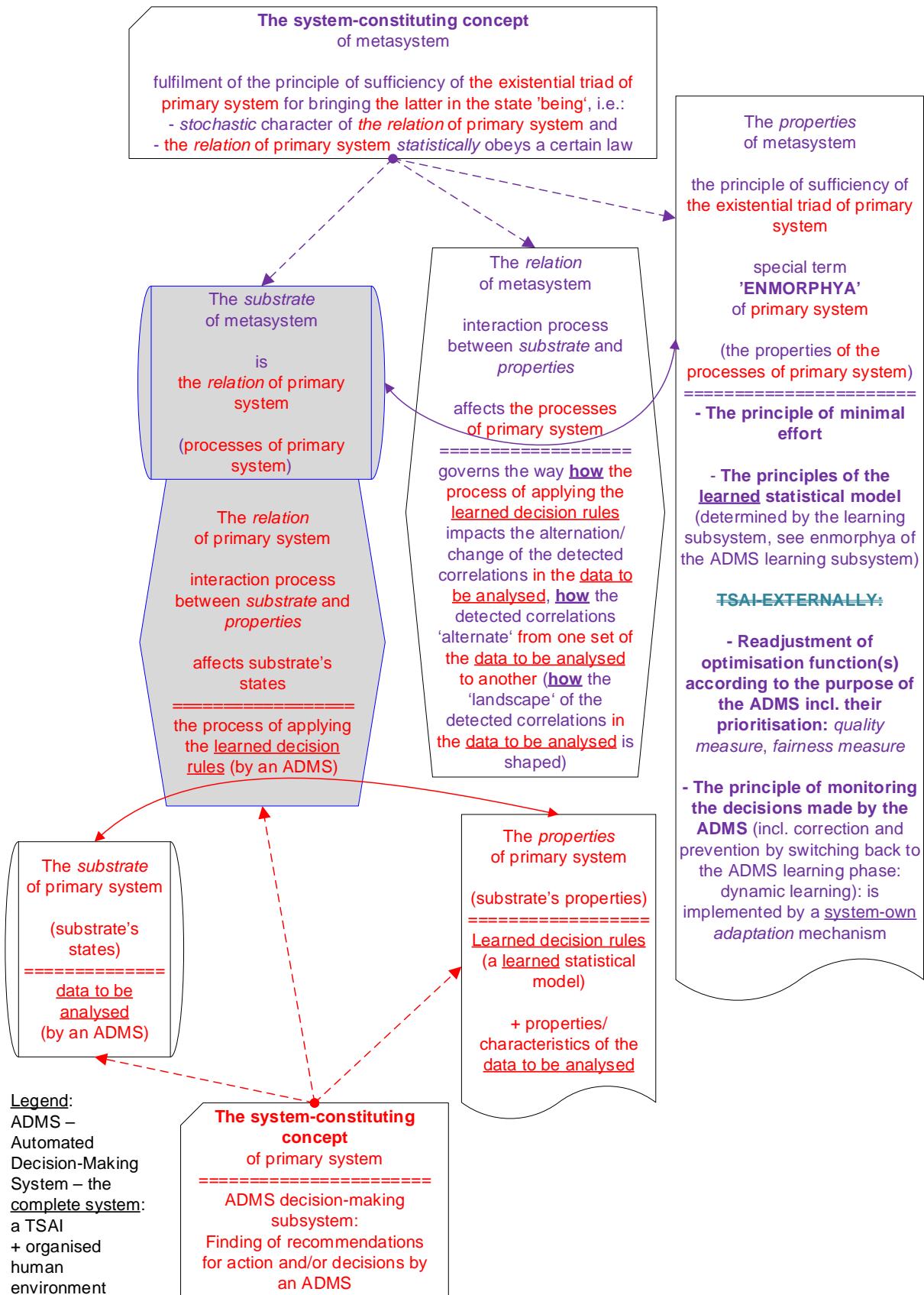


Figure 5: Relationship between the primary system ADMS-DM subsystem (automated decision-making system – the decision-making subsystem) and the corresponding metasytem

If we compare the respective enmorphyas of the relation of the four subsystems of a ‘Weak AI’ on Figure 2 to Figure 5, we can see the systems-theoretical distinctions between these subsystems:

- The TSAI-ML subsystem (Figure 2) is *truly-stochastic* and governed by the *principle of learning economy*; this subsystem is purely technical in nature;
- The ADMS-ML subsystem (Figure 3) is *quasi-stochastic* and thus self-sufficient because this subsystem integrates humans (in the role of ‘data scientist’);
- The TSAI decision-making subsystem (Figure 4) is *deterministic* and governed by the *principles of the learned statistical model*; this subsystem is purely technical in nature;
- The TSAI decision-making subsystem (Figure 5) is *quasi-stochastic* and thus self-sufficient because this subsystem integrates humans (in the role of ‘data scientist’).

This analysis shows, among other things, that purely technical subsystems of a ‘Weak AI’ are either *truly-stochastic* or *deterministic*. The technical AI system – consisting of the *truly-stochastic* TSAI-ML subsystem (the machine learning subsystem) and the *deterministic* TSAI-DM (the decision-making subsystem) – is overall *truly-stochastic*.

This means, among other things, that a ‘Weak AI’ as a purely technical system cannot be an animate system, see [5], CHAPTER II, ch. 1 „Enmorphya of Living Beings“. Thus, the ‘Weak AI’ cannot have the capacity for ‘risk reflection’ either²⁴. Only systems that are capable of asking ‘*essentially necessary questions*’, i.e. only systems possessing *risk reflection*, are capable of creating cognising systems, and thus also capable of semantic cognition²⁵. Applying this insight to a ‘Weak AI’ means:

Statement STM. 1:

A ‘Weak AI’ is fundamentally incapable of a cognition of sense / of a semantic cognition.

This conclusion is in line with the experience gathered to date. Manuela Lenzen writes in [3], ch. 6: “Ultimately, these experiments show that the algorithms do not understand the texts they generate”²⁶. There further she also reproduces the thoughts of John Searle: “Algorithms work with rules that tell you how to move symbols back and forth, they will never grasp meaning. ... Everything that has to do with meaning, Searle says, people interpret into data processing.”²⁷

Janelle Shane also describes similar experiences in several places in [4]: “AI does not really understand the problems to be solved“ (“Introduction”); in ch. 5 therein: “In order for the AI

²⁴ see [5], CHAPTER II, ch. 2.2 „Free Will“, there the new notion of ‘*risk reflection*’; see Glossary below, ‘*risk reflection*’.

²⁵ ‘Cognition of sense’ we synonymously call ‘*semantic cognition*’ (content) as a complementary notion to ‘*syntactic processing of symbols*’ (form).

²⁶ In the original: „Letztlich zeigen diese Experimente, dass die Algorithmen, die Texte, die sie generieren, nicht verstehen.“

²⁷ In the original: „Algorithmen arbeiten mit Regeln, die Ihnen sagen, wie sie Symbole hin- und herschieben sollen, Bedeutung werden sie nie erfassen. ... Alles, was mit Bedeutung zu tun hat, so Searle, interpretieren die Menschen in die Datenverarbeitung hinein.“

to find the right solution, the programmer must ensure that the AI really works on the right problem.”²⁸

An entire ‘Automated Decision-Making System’ (ADMS) is *quasi-stochastic* only thanks to ex-ante and ex-post human intervention. The human (in the role of the ‘data scientist’) must pre-define all semantic goals, continuously monitor their achievement and correctively and preventively readjust ‘optimisation functions’ and ‘hyperparameters’ of a ‘weak AI’.

4 Strong AI

A ‘Strong AI’ is an autonomous technical system that would be universally applicable with the same (or stronger and more extensive) intellectual and creative capabilities as a human being, see [2], Glossary, [3], ch. 1, [4], ch. 2 (“AGI”) or on the Internet.

A ‘Strong AI’ is currently only a vision. A global, intensive social discussion is needed that weighs up the benefits and risks of a ‘Strong AI’. This aspect is discussed in detail in [2], ch. 11.

We ask ourselves here what fundamental prerequisites, from a systems theory perspective, an autonomous technical system shall fulfil so that it would be universally applicable with the same (or stronger and more extensive) intellectual and creative capabilities as a human being.

We have established in [5]²⁹ what distinguishes humans as systems from all other animate systems. The obvious approach to answering the question posed above is to assume that an autonomous technical system should implement the same systems-theoretical principles and essential human-specific attributes of these principles as humans themselves as systems.

Systems-theoretical principles of a system are by definition components of the enmorphya of the relation of this system, see [5], CHAPTER I, ch. 3.2 „Enmorphya“. The principles of the enmorphya of the self-awareness of living beings³⁰ were dealt with in [5], CHAPTER II, ch. 1.2 „Principles of the Enmorphya of Living Beings“: These are the **Principle of Most Choice** and the **Principle of Self-Preservation of System**, see Glossary. However, these two principles are not unique to living beings, but are inherent to any *quasi-stochastic* system³¹.

Thus, we note to ourselves first that a ‘Strong AI’ as an autonomous technical system shall be a *quasi-stochastic* system³², i.e. it shall implement the two principles – of the most choice (of the least resources consumption) and of the self-preservation of system.

This also means, as stated in [5], CHAPTER I, ch. 3.5 „Enmorphya for Quasi-Stochastic Systems“ for any *quasi-stochastic* system, that a ‘Strong AI’ shall necessarily have a system-inherent adaptation mechanism that implements the Principle of Self-Preservation of System. This *adaptation mechanism* includes

²⁸ In the original: „KI versteht die zu lösenden Probleme nicht wirklich“ („Einleitung“); in Kap. 5 darin: „Damit die KI die richtige Lösung finden kann, muss der Programmierer dafür sorgen, dass die KI auch wirklich das richtige Problem bearbeitet.“

²⁹ [5], CHAPTER II, ch. 2.1 „Living and Non-Living Systems“, there in subsection “**The human as a system**”, also in [5], CHAPTER II, ch. 2.2 „Free Will“, there a new notion ‘*risk reflection*’; see Glossary below (‘risk reflection’)

³⁰ The enmorphya of self-awareness is the manifestation, specific to living beings, of the enmorphya of relation.

³¹ The *principle of most choice* is the implementation, specific to living beings, of the general **Principle of Least Resources Consumption**, see [5], CHAPTER II, ch. 1.2 “Principles of the Enmorphya of Living Beings”.

³² Thus, a ‘Strong AI’ shall represent an animate, living system, see [5], CHAPTER II, ch. 2.1 “Living and Non-Living Systems”, and a ‘*will owner*’, see [5], CHAPTER I, ch. 3.5.3 “Society”, [5], CHAPTER II, ch. 2.2 “Free Will” and Glossary.

- *monitoring* of the system state,
- intra-system *correction* (corrective action) with respect to a changing system state, and
- *preventing* a similar ‘sub-optimal’ system state³³ by correcting an appropriate and immanent to that system ‘norm’.

Furthermore, *information metabolism* shall be inherent in a ‘Strong AI’, as in any *quasi-stochastic* system, see Glossary³⁴.

‘Feeling’³⁵ represents a component of *information metabolism*, whereby there can be ‘negative’ and ‘positive’ feelings: A ‘negative feeling’ is felt by a *quasi-stochastic* system when the system is in a ‘suboptimal’ state, i.e. when the stability of its *system-constituting concept* is jeopardised; a ‘positive feeling’ is felt by a *quasi-stochastic* system when the system is in an ‘optimal’ state, i.e. when the stability of its *system-constituting concept* is not jeopardised.

The human-specific (as a species) attributes of the enmorphya of self-awareness are listed in [5], CHAPTER II, ch. 1.3 “Variativity and Attributes of the Enmorphya of Self-Awareness”. These are:

- 1) attribute ‘biological species’ with the value ‘homo sapiens’;
- 2) attribute ‘risk reflection’;
- 3) attribute ‘ethical norms’ as procedural norms valid for a given living system;
- 4) attribute ‘modus’ with possible values ‘ordinary (opportunistic)’ or ‘ontological (ethical)’;
- 5) attribute ‘psychotype’;
- 6) attribute ‘archetype’ in the sense of Jung.

Which of these attributes would be essential for an autonomous technical system with the same intellectual and creative capabilities as a human being to be universally applicable? To answer this question, we apply the method of elimination.

We can directly exclude the attribute ‘*biological species*’ with the value ‘homo sapiens’ because an autonomous technical system cannot be identical to ‘homo sapiens’: Just think of the fact that an autonomous technical system must have a material body other than a human.

The attribute ‘*ethical norms*’, which apply as procedural norms for a human, already exists in the approach for a ‘Weak AI’: This is the *fairness measure*, see ch. 3 above. For a ‘Strong AI’, the respective *fairness measure* is defined by the ‘socialisation’ of the ‘Strong AI’. This attribute is variative. As we noted in [5], CHAPTER I, ch. 3.7.2 „Variativity of Quasi-Stochastic Systems“, the variativity of attributes of enmorphya is an essential distinguishing feature of *quasi-stochastic* from *truly-stochastic* systems.

For a ‘Strong AI’, we note to ourselves here that its enmorphya should include, among other things, at least one specific variative *ethics attribute* that governs the way the ‘Strong AI’ deals with ethical aspects in its communication with its surrounding and its environment.

The attributes ‘*mode*’, ‘*psychotype*’ and ‘*archetype*’ are variative attributes that define the way a human communicates with his surrounding and his environment (in the broadest sense

³³ The ‘sub-optimal’ state of a system is the state that threatens the stability of its *system-constituting concept*.

³⁴ also [5], CHAPTER I, ch. 3.5.3 “Society”, [5], CHAPTER II, ch. 1.1 “Enmorphya of Self-Awareness” and ch. 1.2 “Principles of the Enmorphya of Living Beings”

³⁵ or ‘mood’/‘atmosphere’ for organisations as holistic entities

of the word, i.e. including interaction)³⁶. It is essential that these attributes are not constant but variative.

For a ‘Strong AI’, we note here that its enmorphya should include, among other things, at least one specific variative communication attribute that governs the way the ‘Strong AI’ interacts with its surrounding and its environment.

The attribute ‘*risk reflection*’ represents the main distinctive feature of the free will of humans as a species from the free will of all other animate systems³⁷. Therefore, the presence of this attribute in humans as a species is a prerequisite for human intellectual and creative capabilities.

The human, as the only species, is capable of reflecting on part of his possible (future) states, which include both the world surrounding him and himself, including his own finitude as a system³⁷. The reflection of his own finitude as a system evokes the (predominantly repressed) existential angst, see footnote 36, there subsection “Archetypes” and footnote 37. The existential angst motivates the human to various, among others intellectual and creative activities in order to make the passing of time (always in the direction of the end of the system) unnoticeable in various ways, see [5], CHAPTER VI, ch. 3.4 “The Existential Angst and Adaptation”. It should be noted that angst (and thus also existential angst) is a ‘negative feeling’, see above.

The considerations result in the following chain of dependencies: *Risk reflection* => existential angst (as a ‘negative feeling’)³⁸ => the motivation to make the passing of time unnoticeable => among other things, intellectual and creative activities.

Risk reflection is also a prerequisite for a system to be at all capable of asking ‘*essentially necessary question*’³⁹, i.e. to create cognising systems, and thus also to be capable of semantic cognition⁴⁰.

Some AI researchers introduce an ‘artificial curiosity’ as one of the optimisation functions (in this case - reward or evaluation functions), see [2], ch. 11, [4], ch. 5 (“Curiosity”). This reward function is intended to reward an AI system for finding something more abstract in the ‘world model’ known to the AI system. The intention and hope here is that a ‘Weak AI’ would develop into a ‘Strong AI’ step by step through an ‘artificial curiosity’.

Our considerations above show that ‘artificial curiosity’ – merely as one of the several necessary system-inherent attributes of a ‘Strong AI’ – can only be helpful for a ‘Strong AI’ if this ‘artificial curiosity’ would let a ‘Strong AI’ ask itself ‘*essentially necessary questions*’, i.e. the questions, answers to which are necessary for the implementation of the Principle of Self-Preservation of this ‘Strong AI’³⁹. **Only if this ‘Strong AI’ would reflect its own finitude as a system and feel it as a ‘negative feeling’, it would be self-motivated to ask itself such ‘essentially/existentially necessary questions’.** This means that an ‘artificial curiosity’⁴¹ helpful for a ‘Strong AI’ should have at least these two very specific characteristics.

³⁶ see [5], CHAPTER II, ch. 1.3 “Variativity and Attributes of the Enmorphya of Self-Awareness”

³⁷ see [5], CHAPTER II, ch. 2.2 “Free Will”, STM. 24 and ‘risk reflection’, also Glossary below (‘risk reflection’)

³⁸ Without risk reflection, there can be no existential angst.

³⁹ ‘Essentially necessary questions’ are those questions the answers to which are necessary for the implementation the Principle of Self-Preservation of System. In that sense, ‘*essentially necessary questions*’ can be even called ‘*existentially necessary questions*’, see Glossary.

⁴⁰ Only systems capable of asking ‘*essentially necessary questions*’, i.e. only systems with *risk reflection* are capable of creating cognising systems and thus also capable of cognising sense. ‘Cognition of sense’ we synonymously call ‘*semantic cognition*’ (content) as a complementary notion to ‘*syntactic processing of symbols*’ (form).

⁴¹ A more precise term would be ‘artificial thirst for knowledge’ instead of ‘artificial curiosity’.

For a ‘Strong AI’, we note here that its enmorphya should include the specific attribute ‘*risk reflection*’, which motivates and enables a ‘Strong AI’ to perform intellectual and creative activities, enables it to ask ‘essentially/existentially necessary questions’. In order that *risk reflection* can produce such motivation at all, a ‘Strong AI’ must also experience ‘negative feelings’ with regard to its own finitude as a system.

Now, on the basis of these insights, we can answer the question we posed at the beginning: “What fundamental prerequisites, from a systems theory perspective, an autonomous technical system shall fulfil so that it would be universally applicable with the same (or stronger and more extensive) intellectual and creative capabilities as a human being?”

STM. 2:

In order that a ‘Strong AI’ would be universally applicable as an autonomous technical system with the same (or stronger and more extensive) intellectual and creative capabilities as a human being, the ‘Strong AI’ shall – from a systems theory perspective – implement the following system-inherent principles and corresponding attributes⁴²:

- 1) It should be a *quasi-stochastic* system, i.e. it shall implement both Principles – of Most Choice (of Least Resources Consumption) and of Self-Preservation of System;
- 2) It should necessarily have a system-inherent adaptation mechanism, which implements the Principle of Self-Preservation of System. This *adaptation mechanism* includes
 - *monitoring* of the system state,
 - intra-system *correction* (corrective action) with respect to a changing system state, and
 - *preventing* a similar ‘sub-optimal’ system state⁴³ by correcting an appropriate and immanent to that system ‘norm’;
- 3) Its system-inherent ‘optimisation functions’ and ‘hyperparameters’ shall include, among other things, at least one specific variative ethics attribute that governs the way the ‘Strong AI’ deals with ethical aspects in its communication with its surrounding and environment;
- 4) Its system-inherent ‘optimisation functions’ and ‘hyperparameters’ shall include, among other things, at least one specific variative communication attribute that governs the way the ‘Strong AI’ interacts with its communication with its surrounding and environment;
- 5) Its system-inherent ‘optimisation functions’ and ‘hyperparameters’ shall include, among other things, at least one specific attribute *,risk reflection*⁴⁴ which motivates and enables a ‘Strong AI’ to perform intellectual and creative activities, enables it to ask ‘essentially/existentially necessary questions’;
- 6) Among other things, it should experience ‘negative feelings’ in relation to its own finitude as a system, so that *risk reflection* can first produce such motivation for intellectual and creative activities and for asking ‘essentially/existentially necessary questions’.

⁴² One can also consider principles and corresponding attributes as system-inherent ‘optimisation functions’ and ‘hyperparameters’.

⁴³ The ‘sub-optimal’ state of a system is the state that threatens the stability of its *system-constituting concept*.

⁴⁴ I.e. it shall reflect on a part of possible (future) states, which include both the world surrounding it and itself, including its own finitude as a system.

At the time of publication, the author of these lines had not heard that the need for the above-mentioned system-immanent principles and attributes as prerequisites for the creation of a ‘Strong AI’ was being discussed in the communities dealing with the topic of AI.

5 Enmorphya of ‘Self-Awareness’ of AI: Distinction Criteria between Types of Artificial Intelligence and Perspective

Now we can also establish the systems-theoretical criteria for distinguishing between three types of AI systems that we have considered in this chapter.

Rule-based expert systems are *deterministic*. Their behaviour and properties are apriori predetermined in the corresponding ‘programming manual’.

The purely technical subsystem of a ‘*Weak AI*’ – consisting of the *truly-stochastic* machine learning subsystem and the *deterministic* decision-making subsystem – is as a whole *truly-stochastic*⁴⁵.

The *adaptation mechanism* of a ‘Weak AI’ is not operated technically, but by humans. Its *ethics attribute* – the *fairness measure* – is also specified by humans. I have not yet seen a pendant to the *communication attribute* of a ‘Weak AI’. This attribute could be implemented in a ‘Weak AI’, but the *quality measure* for the corresponding optimisation function would also have to be operated by humans.

A ‘*Strong AI*’, if it is to exist in the future, must be immanently and intrinsically *quasi-stochastic*, so that it would only need the human in the design phase and would be able to act completely autonomously after commissioning. This must include the self-setting of all semantic goals, the self-monitoring of their achievement, and the corrective and preventive self-adjustment of its ‘optimisation functions’ and ‘hyperparameters’. Its system-inherent ‘optimisation functions’ and ‘hyperparameters’ should include the following specific variative attributes: the *ethics attribute*, the *communication attribute* and the ‘*risk reflection*’ among others with regard to its own finitude as a system.

For a ‘Strong AI’, the *adaptation mechanism* and all its attributes should be system-inherent and manage in operation without human intervention.

We have summarised the core differences in the *enmorphya of relation* (in the ‘*enmorphya of self-awareness*’ of AI systems if they were to treat as living beings) between ‘weak AI’ and ‘strong AI’ from a systems theory perspective in the following overview table (differences in **bold**):

⁴⁵ An entire ‘Automated Decision System’ (ADMS) is *quasi-stochastic* thanks to ex-ante and ex-post human intervention. The human (in the role of ‘data scientist’) has to pre-define all semantic goals, continuously monitor their achievement and correctively and preventively readjust ‘optimisation functions’ and ‘hyperparameters’ of a ‘Weak AI’.

AI Type →	<i>Rule-based AI</i> (ch. 2)	<i>'Weak AI'</i> (purely technical subsystem) (ch. 3)	<i>'Strong AI'</i> (if it is to exist in the future) (ch. 4)
Property ↓			
System type	deterministic	truly-stochastic	quasi-stochastic
Main principles of 'optimisation functions'	the contents of the corresponding 'programming manual'	the Principle of Least Resources Consumption (of Most Choice)	- the Principle of Least Resources Consumption (of Most Choice) - the Principle of Self-Preservation of System
Adaptation mechanism	not applicable	is adjusted by the <u>system-external</u> human being	is system-inherent and self-adjusting
Ethics attribute	not applicable	The <i>fairness measure</i> is adjusted by the <u>system-external</u> human being.	variative attribute of system-inherent 'optimisation functions' and 'hyperparameters'
Communication attribute	not applicable	We have not seen it yet, but it would be applicable and implementable; The <i>quality measure</i> for the corresponding optimisation function would have to be adjusted by the <u>system-external</u> human.	variative attribute of system-inherent 'optimisation functions' and 'hyperparameters'
Risk reflection ⁴⁶	not applicable	not applicable	an attribute of system-inherent 'optimisation functions' and 'hyperparameters'; motivates and enables intellectual and creative activities as well as asking of 'essentially/existentially necessary questions' (see

⁴⁶ I.e. it shall reflect on a part of possible (future) states, which include both the world surrounding it and itself, including its own finitude as a system.

AI Type →	<i>Rule-based AI</i> (ch. 2)	<i>‘Weak AI’</i> (purely technical subsystem) (ch. 3)	<i>‘Strong AI’</i> (if it is to exist in the future) (ch. 4)
Property ↓			Glossary)
Feeling ‘negative feelings’ about their own finitude as a system ⁴⁷	not applicable	not applicable	a result of the function of the <i>monitoring mechanism</i> within the <i>adaptation mechanism</i>
Capability to cognition of sense ⁴⁸	none	none	fundamentally possible due to <i>risk reflection</i>, which motivates and enables the system to ask ‘<i>essentially/existentially necessary questions</i>’

This overview shows that the core differences between ‘weak AI’ and ‘strong AI’ lie – from a systems theory perspective – in the *enmorphotype*⁴⁹ of the specific technology implementing a specific AI.

Our systems-theoretical consideration shows, among other things, that ‘Weak AI’ and ‘Strong AI’ are significantly different from each other. It also shows that the creation of a ‘Strong AI’ may have a very long way to go, if it should ever become possible at all. The global ethical aspect of whether humanity even wants such a creation requires a global intensive societal discussion about benefits and risks.

⁴⁷ As a result of the function of the *monitoring mechanism* within the *adaptation mechanism*.

⁴⁸ ‘Cognition of sense’ we synonymously call ‘*semantic cognition*’ (content) as a complementary notion to ‘*syntactic processing of symbols*’ (form).

⁴⁹ see [5], CHAPTER II, ch. 1.4 “*Enmorphotype*”

6 Glossary

Term	Definition
Basic notions of systems theory by A. Uemov [1], necessary for reading this work	
system	<p>Any given entity on which a <i>relation</i>, possessing an arbitrarily taken certain <i>property</i>, is implemented.</p> <p>Or equivalently:</p> <p>any given entity on which some <i>properties</i>, being in an arbitrarily taken certain <i>relation</i>, are implemented.</p>
system-constituting concept ⁵⁰	An a priori given system-constituting <i>property</i> or <i>relation</i> ; dependent on this, a system-constituting concept is an <i>attributive</i> or <i>relational</i> one, resp.
structural factor ⁵¹	<p>A set of properties and relations that suffices the given system-constituting concept.</p> <p>A structural factor can be a relational one (in the case of the attributive concept) and an attributive one (in the case of the relational concept).</p>
system substrate ⁵²	A carrier of a relational or attributive structure.
Other basic notions necessary for reading this work	
existential triad	<p>A set of {<i>substrate</i>, <i>property</i>, <i>relation</i>} that is necessary for creating a system based on this set.</p> <p>An existential triad is sufficient for the creation of a system with its corresponding <i>system-constituting concept</i>, if the ‘relation’ in this triad</p> <ul style="list-style-type: none"> - is fundamentally <i>stochastic</i>, and - <i>statistically</i> obeys a certain law (in the general case – the PLR – the Principle of Least Resources Consumption). <p>The evolution of this system follows the character of the ‘relation’ in the existential triad.</p>
universal existential pentad	<p>A form necessary and sufficient to describe the abstract structure of <u>any</u> system (and thus <u>any</u> observable entity) regardless of the content and purpose of that system and the principles governing that system.</p> <p>The universal existential pentad is the whole schema itself, shown in Figure 1 in [5], Part A, Chapter I ch. 3.2 “Enmorphya”, i.e. all five elements of the schema and the relationships between these elements, i.e.:</p> <ul style="list-style-type: none"> - the <i>substrate</i> of the primary system, - the <i>properties</i> of the primary system, - the <i>relation</i> of the primary system = the <i>substrate</i> of the metasystem, - the <i>properties</i> of the metasystem (<i>enmorphya</i> of the relation), and - the <i>relation</i> of the metasystem.

⁵⁰ the original term by A. Uemov: ‘системообразующий концепт’

⁵¹ the original term by A. Uemov: ‘структурный фактор’

⁵² the original term by A. Uemov: ‘субстрат системы’

Term	Definition
	<p>The <i>existential pentad</i> is <u>universal</u> and <u>complete</u>.</p>
information	<p>A change in the degree of indeterminacy</p>
information metabolism	<p>The existentially necessary reception and processing of signals from the environment by the system and the system's response to these signals.</p> <p>Information metabolism is inherent not only in a person, but also in any <i>quasi-stochastic</i> system because its <i>adaptation</i> mechanism cannot work without exchanging and processing signals with the environment.</p> <p>The concept of 'information metabolism' was introduced by Antoni Kępiński as a parallel to the energy metabolism of the body (Antoni Kępiński <i>Psychopatologia nerwic (Psychopathology of Neuroses)</i>, 1972).</p>
essentially necessary questions	<p>Questions the <u>answers</u> to which are necessary for the implementation the Principle of Self-Preservation of System.</p> <p>In that sense, '<i>essentially necessary questions</i>' can be even called '<i>existentially necessary questions</i>'.</p>
adaptation	<p>An adjustment of an intra-system 'norm' (changing it, abolishing it, creating a new one) as a result of the effect of <i>feedback</i></p> <p>The <i>adaptation</i> mechanism comprises the mechanisms for</p> <ul style="list-style-type: none"> - <i>monitoring</i> of the system state (which also depends on environmental conditions), - intra-system <i>correction</i> (corrective action) with respect to a changing system state, and - <i>preventing</i> a similar 'sub-optimal' system state by correcting an appropriate, immanent to that system 'norm'. <p>These mechanisms are immanent to the system.</p> <p>The 'sub-optimal' state of a system is the state that threatens the stability of its <i>system-constituting concept</i>.</p> <p>The combination of <i>monitoring</i> and <i>correction</i> mechanisms is often referred to as a <i>feedback</i> mechanism.</p> <p>For <i>quasi-stochastic</i> systems, all three of these mechanisms exist and must be active.</p> <p>For <i>truly-stochastic</i> systems, which have no <i>long-term memory</i>, the <i>prevention</i> mechanism cannot function, as the <i>long-term memory</i> is necessary to maintain the intra-system 'normative base'. Therefore, the <i>adaptation</i> mechanism for <i>truly-stochastic</i> systems is equivalent to the <i>feedback</i> mechanism (<i>monitoring</i> and <i>correction</i> only).</p>
resource (of a system)	<p>The product 'number of steps on the way from state A to state B' by 'number of alternative solutions/opportunities at each such step'.</p> <p>The resource of the system can be abstractly represented as the product of two categorially complementary terms:</p> <p style="text-align: center;">'resource' = 'action' * 'choice',</p> <p>see details in [5], CHAPTER VII, ch. 2.3.2 "Complementary Terms as Resource".</p>

Term	Definition
	<p>The specific implementation of ‘steps on the way from state A to state B’ and ‘alternative solutions/opportunities at each such step’, i.e. the specific implementation of ‘action’ and ‘choice’, is specific in each system and must be defined for each system separately⁵³.</p> <p>For example, for physical systems the ‘resource’ is the number of action quanta necessary to transition the system to another given macroscopic state⁵⁴; for communication (including the communication function of language) – the number of positions in the message (text) * the number of different signs (for example, letters and punctuation marks) necessary to convey the given content; for educational – and for any other social process – the number of individual (learning) topics * the number of alternative (didactic) methods to be considered and applied, respectively, for the achievement of a given (learning) objective.</p>
the Principle of Least Resources Consumption (PLR)	<p>The principle of dynamics of development of any system that consists in the fact that a system at transition from state A to state B implements in statistical average such a way of transition from A to B, at which the ‘resource’ of the system is consumed at the least.</p> <p>PLR is a universal relation-control-information (i.e. is an integral part of <i>enmorphya</i> of relation) and governs the process of <i>interaction</i> between the <i>substrate</i> and the <i>structural factor</i> of <u>any</u> system – physical, social, communicative, etc. – which is based on a <i>stochastic</i> process.</p> <p>In particular, the PLR governs the process of interaction between matter and information in nature in the form of the principle of most entropy which is equivalent to the principle of least action, cf. [5], CHAPTER VII, ch. 2.1.5 “The Principle of Least Resources Consumption” and 2.3.2 “Complementary Terms as Resource”.</p>
the Principle of Self-Preservation of System (PSP)	<p>The principle of stabilisation of any system, which consists in the fact that the <u>deviation</u> of the system from obeying the Principle of Least Resources Consumption is limited by the fact that the system-constituting concept of this system remains stable, retains.</p> <p>The Principle of Self-Preservation of System is valid for <u>any</u> system, i.e. it is a universal part of their <i>enmorphya</i>. For <i>truly-stochastic</i> systems, it is done automatically due to their ‘being Markovian’, which in itself brings the stochastically ‘out of line’ systems back to the path of maximum entropy.</p> <p>For <i>quasi-stochastic</i> systems, there is no such automatism. Its absence shall therefore be compensated for by the system’s explicit, inherent mechanisms to help preserve the system. Such (system-immanent) mechanisms are implemented through an <u>adaptation mechanism within the system itself</u>.</p>
the Principle of Most Choice	<p>The principle of minimizing the restrictive factors on the opportunities of making decisions, the principle of maximizing the freedom of choice.</p> <p>It is the Principle of Most Choice as one of the characteristics of the self-awareness of living beings that leads to their flexibility and <u>adaptability</u> to various conditions of existence.</p>

⁵³ The number of ‘steps on the way from state A to state B’ must be > 0 , and the number of ‘alternative solutions/opportunities at each such step’ must be > 1 . The reason for this is that nature must spend more than zero resources to create an observable state. For this, nature ‘must’ make at least 1 ‘step to another state’ and ‘alternative solutions at each such step’ cannot be deterministic and therefore the number of alternatives must be > 1 ; see [5], CHAPTER VII, ch. 2.1.3, 2.1.4, 2.3.2 for further details.

⁵⁴ i.e. the physical quantity ‘action’ ($\text{kg}\cdot\text{m}^2\cdot\text{s}^{-1}/\text{h}$) (the Planck constant is the value of action quantum)

Term	Definition
enmorphya ⁵⁵ of <i>sth.</i>	<p>A particular term for the notion ‘control-information-of-<i>sth.</i>’, e.g. ‘enmorphya of <i>relation</i>’</p> <p>The distinguishing mark between the notions ‘information’ and ‘enmorphya’ consists in the following: ‘information’ interacts with <u>material <i>substrate</i></u>, whereas ‘enmorphya’ interacts with <u>the relation and process</u> between this ‘information’ and this material <i>substrate</i>.</p>
stochastic process	A process whose every next state occurs with any probability other than 0 and 1
stochastic system	A system whose <i>structural factor</i> is based on a <i>stochastic process</i>
deterministic process	<p>A process whose every next state is <u>unambiguously defined</u> by its present state, i.e. every next state comes with probability 1.</p> <p>This means that each previous state of the process can also be unambiguously calculated from its present state.</p> <p>If the next process state comes with probability 0 then the process has stopped and no longer exists; it also falls within the definition of deterministic process.</p>
deterministic system	A system whose <i>structural factor</i> is based on a <i>deterministic process</i>
Markov property (of a <i>stochastic process</i>)	<p>Every next state of the Markov <i>stochastic process</i> implementing regular Markov chains probabilistically depends <u>solely</u> on its current state and is independent of its previous states.</p> <p>This property can also be expressed in the following way: the past of the <i>truly-stochastic</i> (i.e. Markovian) systems affects their future exclusively through their present.</p>
truly-stochastic process	<p>A <i>stochastic process</i> possessing the <i>Markov property</i></p> <p><u>The ‘true stochasticity’ is the absence of immediate (direct) memory</u> of previous states: the subsequent state probabilistically depends only on the current state.</p> <p>The <i>enmorphya</i> of relation is <u>non-variable</u> (always the principle of least action without variable characteristics).</p>
quasi-stochastic process	<p>A <i>stochastic process</i> that has <u>no Markov property</u></p> <p><u>Quasi-stochastic</u> systems must possess <u>immediate (direct) long-term memory</u> of previous states.</p> <p>The <i>enmorphya</i> of relation is <u>variable</u> (always the <i>Principle of Least Resources Consumption</i> with variable characteristics and the <i>Principle of Self-Preservation of System</i> with an <i>adaptation</i> mechanism).</p> <p>N.B.: <i>quasi-stochastic</i> processes are <u>not deterministic</u>.</p>
will owner	<p><u>Any quasi-stochastic</u> system, i.e. a stochastic system with <i>freedom of choice</i>, which takes into account all its previous experience and has an <i>adaptation</i> mechanism</p> <p>In other words, a <i>will owner</i> is an <i>adaptive</i> system with <i>freedom of</i></p>

⁵⁵ The term ‘enmorphya (enmorfia, enmorphy)’ is constructed on the basis of Greek: ἐνμορφία (ἐν-μορφή-α => (bringing) in-form, (приведение) в-форму)

Term	Definition
	<i>choice.</i>
socium	A social entity/unit, a group of <i>will owners</i> , a socially connected system of interacting <i>will owners</i> , a society of any size held together by any internal relationships
categorial complementarities	Let there exist a confined population (set) of terms comprising more than one term. Terms out of the population are called <i>categorially complementary</i> to each other if: <ol style="list-style-type: none"> 1) These terms can exist exclusively jointly, in concert, i.e. the existence of a term necessarily causes the existence of all other terms of the population, and 2) A term out of the population cannot be defined by using any subset of other terms of the population.
attributive opposites	Let there exist a confined population (set) of properties comprising more than one property. Properties out of the population are called <i>attributive opposites</i> if each item of the population represents merely a <u>specific extreme value of one and the same attribute</u> , and, hence, can be defined by using another item of the population. <p>Distinguishing between <i>attributive opposites</i> (e.g. {high, low}) and <i>categorial complementarities</i> (e.g. {form, content}), let it be said that attributive opposites are basically not categorial complementarities because each item of an attributive pair can be defined by using another member of the pair. For example, the attribute ‘size’ can take extreme values {big, small}; these values can be expressed by each other.</p> <p>Attributive opposites always describe properties/qualities, i.e. <u>values of an attribute</u>, but never – terms. Thereby, changing the value of this attribute at the transition from one to another extreme occurs without ‘jumps’, i.e. without a change of symmetry degree (without ‘second-order phase transitions’). Attributive opposites often imply the presence of an etalon, i.e. a ‘norm’, which the estimation of the value of the respective attribute relates to (e.g. {expensive, cheap}, {good, evil}).</p> <p>Attributive opposites almost always are reflected in language by antonymous pairs, whereas <u>categorial complementarities are by no means always representable by them.</u></p>
time	Distinguishability of the microstates of nature from each other IS the course of time (i.e. time itself). <p>Therefore, time is discrete.</p> <p>Distinguishability of states is a necessary prerequisite for their observability, i.e. for their being. That is why being and time are bijectively connected with each other. See [5], CHAPTER VII, ch. 1.3 “Time Microstructure”.</p>
past	Recorded/documented set of states (events) that have occurred. <p>Therefore, the past is deterministic, see [5], CHAPTER VII.</p>
the present	Decision-making on choosing the next state from a variety of possible states. <p>The present turns a probabilistic future into the deterministic past. It is this complementarity of the probabilistic future and the deterministic past that causes the <i>irreversibility</i> of time, see [5], CHAPTER VII.</p>

Term	Definition
instant	<p>A theoretical notion describing an ‘intermediate state’ that cannot be realised in nature.</p> <p>In such an ‘intermediate state’, the possibility of choice already exists but the resolution of this alternative does not exist yet. Since time is discrete, there cannot be any ‘intermediate states’ of entities.</p> <p>This definition makes the ‘instant’, and with it the present, a relative, but not an absolute notion.</p>
future	<p>A variety of possible states.</p> <p>Therefore, the future is probabilistic, see [5], CHAPTER VII.</p>
memory	<p>The property of storing information (both rational and emotional, if applicable to a given system) for a period of time beyond a given state (<i>instant</i>, situation) of the system, so that this stored information can directly affect <u>more than one</u> subsequent state (situation) of that system</p> <p>Such memory can also be called ‘<i>long-term memory</i>’.</p> <p><i>The long-term memory</i> is a necessary attribute of the <i>quasi-stochastic process</i>.</p> <p>In this context, ‘<i>short-term memory</i>’ is the property of storing information (both rational and emotional, if applicable to a given system) for a period of time not exceeding the given state (<i>instant</i>, situation) of the system, so that this stored information can directly affect no more than one – the next – subsequent state (situation) of that system.</p> <p><i>Short-term memory</i> realises the <i>Markov property</i> and is a necessary attribute of the <i>truly-stochastic process</i>.</p>
history	<p>The sequence of phases in the development of the <i>quasi-stochastic</i> system, i.e. of the <i>will owner</i> to whom this ‘history’ pertains.</p> <p>The full history of the <i>will owner</i> includes the full cycle of development of the corresponding <i>quasi-stochastic</i> system from its emergence to its self-destruction. This complete cycle of development exists for any <i>quasi-stochastic</i> system.</p>
space	<p>A discrete <i>substrate</i> needed for <i>distinguishing</i> between <u>material</u> entities, see [5], CHAPTER VII, ch. 3 “Space Microstructure”.</p>
enmorphotype (of a living being)	<p>A set of all attributes of the ‘<i>enmorphya</i> of self-awareness’ of a living being which interacts with both his/her genotype and phenotype.</p>
free will	<p>Free will is the freedom of choice, which is non-deterministic, but does not represent a Markov process and takes into account at least all previous experience of the system.</p> <p>I.e. it is a certain freedom of choice, a possibility of local deviation of <i>quasi-stochastic</i> process from following the Principle of Least Resources Consumption.</p> <p>The decision-making process.</p>
risk reflection (by human beings only) (uncertainty of the possible (of the future))	<p>Inclusion in decision-making, i.e. in the freedom of <u>human</u> choice, of a self-reflection of possible future states that include both the world surrounding the human and the human itself, including its own finitude as a system.</p>

7 References

- [1] Авенир Иванович Уёмов *Системные аспекты философского знания*, Одесса, 2000⁵⁶
- [2] Katharina Zweig *Ein Algorithmus hat kein Taktgefühl*, Heyne, 2019, ISBN 978-3-453-20730-1
- [3] Manuela Lenzen *Künstliche Intelligenz*, C.H.Beck, 2020, ISBN 978-3-406-75124-0
- [4] Janelle Shane *You Look Like a Thing and I Love You*, 2019 Janelle Shane, ISBN 978-1-47226-899-0
- [5] Igor Furgel *Being and Systemacy*, Westarp BookOnDemand, 1. Auflage 2022, ISBN: 978-3-96004-131-3

8 Acknowledgements

I would like to express my deep gratitude to my wife Irina for our extremely useful and interesting discussions on certain aspects of this topic. I would also like to express my deep gratitude to my university professor of philosophy Avenir I. Uemov for his invaluable participation in shaping my style of interaction with the world.

⁵⁶ Avenir Ivanovich Uemov *Systems Aspects of Philosophical Knowledge*, Odessa, 2000