Is Your AI QUALITY Good Enough To Save Humanity?

Artificial Narrow Intelligence (ANI)



Stage-1

Machine Learning

Specialises in one area and solves one problem



Artificial General Intelligence (AGI)



Stage-2

Machine Intelligence

Refers to a computer that is as smart as a human across the board

Artificial Super Intelligence (ASI)



Stage-3

Machine Consciousness

An intellect that is much smarter than the best human brains in practically every field

https://www.mygreatlearning.com/blog/what-is-artificial-intelligence/

TOM GILB BCS 3 July 2024

45 minutes talk and 15 minutes Q&A Recorded (1 MIN./SLIPE)

tom@Gilb.com

Gilb International (ret.) Norway

These slides folder (+ appendix) is-your-ai-quality-good-enough-to-save-humanity SLIPES FOLPER : https://www.dropbox.com/scl/fo/ ulfisxhd5cyrp7iij2i98/AAYYSATm4Yt40X0t8lEJdpQ? rlkey=3n4nn6077nig6gvo5giuu0em8&dl=0



SLIDES FOLDER DOI to slides pdf: <u>https://doi.org/10.13140/RG.2.2.10935.12960</u>



- 1. Current ways to 'understand' AI systems 2. Why these current AI Models are defective. **3. Multidimensional AI Models for Qualities, And Costs**
- 4. The AI Stakeholder Model
- 5. Understanding AI Strategies in light of Multiple Quality and Cost Attributes 6. Teaching The Multi-Al model to students and research students as a tools
- for researching AI developments.
- 7. The Quality distinctions between Large Language Model (LLM) AI and the next AI generation (AGI, Artificial General Intelligence) 8. The Penta Model as a basic high level view of any Al system. 9. The Principles of Al Understanding, outside the blackbox.

Outline: Al QUALITY

1:41 / 51:10



Corners uh on safety



Al Safety According To Tesla, The Model 3 Is The Safest Car On The Planet

OCT 9, 2018 / BY JAY TRAUGOTT IN TESTING 🤊

And the number two and three spots go to ...

Last month it was announced the new Tesla Model 3 earned a perfect five-star crash test rating in every category from the National Highway Transportation Safety Administration. Most impressive. And now Tesla says its new mid-size sedan has the lowest probability of causing injury for any vehicle ever tested by the US government agency.

"Not only has Model 3 achieved a perfect 5-star safety rating in every category and sub-category, but NHTSA's tests also show that it has the lowest probability of injury of all cars the safety agency has ever tested," Tesla wrote in a blog post. This same post also explains why the Model 3 is so safe.



Are you an Al Engineer, yet?

1. How many of you know how to quantify (define a metric for) Al System Transparency, and can show me you have done it, in writing?

3. How many of you know how to quantify (define a metric for) Al System 'Security', and can show me you have done it, in writing?



Some insights

Into A

methods



https://altern.ai/topic/resources



1. Current (BAD) ways to 'understand' Al systems

Design ('How', inside the black box)

- **PESIGNS:** We use Big Pata, LLM (Large Language Models), 'Pata * **Reflectivity**'
- Functions ('What it does', inside black box) "It does poems and art" *
- Values-Qualities ('How Well', a result of 'Design'), like Security, Usability, Transparency, Traceability.
- "It is easy to use, but not very trustworthy, and difficult to * understand, why it gives such answers."
- *

*

*

Usually described in imprecise terms, non-numeric, unmeasurable 'easy use', 'fairly secure') (like

Resources ('Costs', a result of 'Design', and 'Operate')



*

"It is free to play with, expensive to build, and answers quickly.

A large language model: LLM is a language model consisting of a neural network with many parameters, trained on large quantities of unlabelled text using self-supervised learning.

LLMs emerged around 2018





2.1 Why these current AI 'Models' (ways of describing an AI system) are defective (bad).

- Pesign ('How', inside the black box)
 - Pesign is secret, hidden, non-transparent, un-intelligible, complex, changing (updates), expanding automatically
- * Functions ('What', inside black box, part of 'Scope')
 - * A result of the Design, same problems as above (un-intelligible)
 - * What kind of poems? What kind of art? (Highly ambiguous)
- * Values-Qualities ('How Well', a result of 'Design'), like Security, Usability, Transparency......
 - * a result of Unintelligible Design. Unstable, unpredictable
 - susually described in imprecise terms, non-numeric, unmeasurable (like 'easy use', 'fairly secure')
 - * there are very many (20+) interesting qualities all at once
- * Resources ('Costs', a result of 'Design')
 - Same problem as Qualities, for same reasons: just nice words (but, <u>https://ai-benchmark.com/</u>)



2.2 Why these current AI 'Models' (ways of describing an AI system) are defective.

- * The technical structure of the Al system is inherently unstable, because...
- * The <u>conscious</u> design (red triangle) (the program and algorithms, by people)
 - * can be changed by <u>unidentified</u> people
 - * for <u>unidentified</u> reasons,
 - and even introduce <u>unintended</u> side-effects (bugs),
 - * and receive insufficient quality assurance and testing (until proven otherwise).



2.3 Why these current AI 'Models' (ways of describing an AI system) are defective.

*

The technical structure of the Al system is inherently unstable, because...

* the Bot-Design (Effectively = changes to the Human Design, via Data),

* and enabled by the (ongoing) human redesign)

is done automatically, instantaneously, and by rules dictated by 'Casual Pata'

= anything on internet, at any time, or generated by itself (AI), or fed by Hackers, including Fake News)

Casual Pata: data that just happens'



2.4 Why these current AI 'Models' (ways of describing an AI system) are defective.

- * The technical structure of the Al system is inherently unstable, because...
- * THE CONCLUSION
 - * You cannot know the 'design' + casual data - that gave you some answer'
 - * So knowing the 'design' (inside the Al Black box) is not a good way to understand the Al system
 - * in order to decide if is truthful, accurate, false, fake, useful.





Mystical Complexity in the Al System

Understanding **Critical Al** Properties By guantification and measurement, of, for example, **TRANSPARENCY**





Sources in presenter notes

10 Non-Existent 'facts' per report

Factuality

Unintelligible **Complexity inside** the AI Black box

Technoscopes



Mystical Complexity in the AI System

42 Intentional **Fake news** per report

Misleadingness



https://app.simplified.com/tools/text-to-image/my-creations

AI TOOL PROMPT: "a black box with chaos inside and 5 sets of 3 peepholes on outside of it one in each side, 3d rotating"



DALL-E. https://labs.openai.com/e/NfHUMF9htol5j847yGwGQc4i/ boxHJS683d7Mlhs91vsu0Giz





Full interview: "Godfather of artificial intelligence" talks impact and potential of AI



CBS Mornings
2.5M subscribers



ribe





493K views 3 weeks ago #artificialintelligence #tech

https://www.youtube.com/watch?v=qpoR0378qRY

Geoffrey Hinton is considered a <u>godfather of artificial intelligence</u>, having championed machine learning decades before it became mainstream. As chatbots like ChatGPT bring his work to widespread attention, we spoke to Hinton about the past, present and future of AI. CBS Saturday Morning's Brook Silva-Braga interviewed him at the Vector Institute in Toronto on March 1, 2023.

Hinton is mainly interested in Studying how the brain works.

So we need to note that He says the Al we are developing is 'very different from our intelligence'

A Share

...

How different, in what dimensions or attributes.....?

This is the question I am trying to find a way to answer, in this talk. TG



3.1 Multidimensional AI Models for Qualities, And Costs

- * So here is my advice... (. how to understand 'A.I.')
 - * (And it is good general advice for complex systems)**
 - * Measure qualities and other stakeholder values
 - * Use a set of Al 'qualities', to 'understand' an Al system, in terms of what it means for stakeholders (people)
 - * The total of this qualities set defines Al 'Effectiveness'
 - * Use a set of Al-system costs to understand 'economics'
 - * The total set of these costs helps us understand to Alresources needed.
 - * The Values/Costs helps us understand the 'efficiency' of the Al-system







3.2 Multidimensional AI Models for Qualities, And Costs

- * Here are some examples of AI Quality sets
 - * IEEE Ethics Model
 - * 1. Human Benefit
 - * 2. Responsibility
 - * 3. Transparency
 - * 4. Education and Awareness

https://standards.ieee.org/wp-content/uploads/import/documents/other/ead_general_principles.pdf sk

* Al for The White House



XAI Talk 2019 Slides



Develop new standards that describe measurable, testable levels of transparency, so that systems can be objectively assessed and levels of compliance determined. For designers, such standards will provide a guide for self-assessing transparency during development and suggest mechanisms for improving transparency. (The mechanisms by which transparency is provided will vary significantly, for instance (1) for users of care or domestic robots a why-did-you-do-that button which, when pressed, causes the robot to explain the action it just took, (2) for validation or certification agencies the algorithms underlying the AI/AS and how they have been verified, (3) for accident investigators, secure storage of sensor and internal state data, comparable to a flight data recorder or black box.)

I Quantified, in later slides

PREPARING FOR THE PUTURE OF ARTIFICIAL INTELLIGENCE of Schener and Technology Exe mitter on Technolog Dahebar 2020

© 2017 In maliane Bosmess Medimes Coussea its

Proposal of Discussion toward Formulation of AI R&D Guideline

Referring OECD guidelines governing privacy, security, and so on, it is necessary to begin discussions and considerations toward formulating an international guideline consisting of principles governing R&D of AI to be networked ("AI R&D Guideline") as framework taken into account of in R&D of AI to be networked.

Proposed Principles in "AI R&D Guideline"

- Principle of Transparency
- Ensuring the abilities to explain and verify the behaviors of the AI network system
- Principle of User Assistance

Oct. 2016 : "Preparing for the future of AI" from USA

- Giving consideration so that the AI network system can assist users and appropriately provide users with opportunities to take choices
- . Principle of Controllability
- Ensuring controllability of the Al network system by humans
- . Principle of Security
- Ensuring the robustness and dependability of the AI network system Principle of Safety
- Giving consideration so that the AI network system will not cause danger to the lives/bodies of users and third parties **Principle of Privacy**
- Giving consideration so that the AI network system will not infringe the privacy of users and third parties . Principle of Ethics
- Respecting human dignity and individuals' autonomy in conducting research and development of AI to be networked Principle of Accountability
- Accomplishing accountability to related stakeholders such as users by researchers/developers of AI to be networked



IBM

17

Oct. 2016 : "Preparing for the future of Al" from USA



Distributed materia **Proposal of Discussion toward Formulation of Al R&D Guideline**

Referring OECD guidelines governing privacy, security, and so on, it is necessary to begin discussions and considerations toward formulating an international guideline consisting of principles governing R&D of AI to be networked ("AI R&D Guideline") as framework taken into account of in R&D of AI to be networked.

Proposed Principles in "Al R&D Guideline"

1. Principle of Transparency

Ensuring the abilities to explain and verify the behaviors of the AI network system

2. Principle of User Assistance

Giving consideration so that the AI network system can assist users and appropriately provide users with opportunities to make choices

3. Principle of Controllability

Ensuring controllability of the Al network system by humans

4. Principle of Security

Ensuring the robustness and dependability of the Al network system

- 5. Principle of Safety
- 6. Principle of Privacy

Giving consideration so that the AI network system will not infringe the privacy of users and third parties

7. Principle of Ethics

Respecting human dignity and individuals' autonomy in conducting research and development of AI to be networked

8. Principle of Accountability



Giving consideration so that the AI network system will not cause danger to the lives/bodies of users and third parties

Accomplishing accountability to related stakeholders such as users by researchers/developers of AI to be networked





"SERIOUSLY guys, don't try to bullshit a bullshitter!"

- * "Ensuring the abilities to explain and verify the Al Network System"
 - * 7 ambiguous words
 - * "Ensuring":
 - improving the probability, (2 ambiguities here)
 - * guaranteeing,
 - * making sure it is perfectly explained,
 - * making sure it is really perfectly effective
 - * And 500 more



Proposal of Discussion toward Formulation of AI R&D Guideline

Referring OECD guidelines governing privacy, security, and so on, it is necessary to begin discussions and considerations toward formulating an international guideline consisting of principles governing R&D of AI to be networked ("AI R&D Guideline") framework taken into account of in R&D of AI to be networked

Proposed Principles in "AI R&D Guideline"

- Principle of Transparency
- Ensuring the abilities to explain and verify the behaviors of the AI network system
- Principle of User Assistance

Giving consideration so that the AI network system can assist users and appropriately provide users with opportunities to ake choices

- . Principle of Controllability
- Ensuring controllability of the Al network system by humans
- Principle of Security
- Ensuring the robustness and dependability of the AI network system Principle of Safety
- Giving consideration so that the AI network system will not cause danger to the lives/bodies of users and third parties Principle of Privacy
- Giving consideration so that the AI network system will not infringe the privacy of users and third parties Principle of Ethics
- Respecting human dignity and individuals' autonomy in conducting research and development of AI to be networked . Principle of Accountability
- Accomplishing accountability to related stakeholders such as users by researchers/developers of AI to be networked

PREPARING FOR THE FUTURI OF ARTIFICIAL INTELLIGENCE

Science and Technology (

nittee on Techno

October 2016



17

... we need to design and report on the <u>metrics</u> ("You can't regulate what you don't understand.")

https://www.oreilly.com/content/you-cant-regulate-what-you-dont-understand-2/

But whose human values? Those of the benevolent idealists that most AI critics aspire to be? Those of a public company bound to put shareholder value ahead of customers, suppliers, and society as a whole? Those of criminals or rogue states bent on causing harm to others? Those of someone well meaning who, like Aladdin, expresses an ill-considered wish to an allpowerful Al genie?

There is no simple way to solve the alignment problem. But alignment will be impossible without robust institutions for disclosure and auditing. If we want prosocial outcomes, we need to design and report on the metrics that explicitly aim for those outcomes and measure the extent to which they have been achieved. That is a crucial first step, and we should take it immediately. These systems are still very much under human control. For now, at least, they do what they are told, and when the results don't match expectations, their training is quickly improved. What we need to know tim@oreily.com is what they are being told.



https://www.oreilly.com/tim/ tim@oreilly.com

Thanks for this Pawel Nowak



3.3 Multidimensional AI Models for Qualities (quantified), And Costs Here are some examples of AI Quality sets

* IEEE Ethics

*

* 1. Human Benefit

* 2. Responsibility

*** 3. Transparency**

* 4. Education and Awareness

https://standards.ieee.org/wp-content/uploads/import/ documents/other/ead_general_principles.pdf Candidate Recommendation

Develop new standards that describe measurable, testable levels of transparency, so that systems can be objectively assessed and levels of compliance determined.

For designers, such standards will provide a guide for self-assessing transparency during development and suggest mechanisms for improving transparency. (The mechanisms by which transparency is provided will vary significantly, for instance (1) for users

of care or domestic robots a why-did-you-do-that button which, when pressed, causes the robot to explain the action attens) to explain the action attens) to explain the action attens) to explain the action attens and explain the action attens attended to be action at the algorithms underlying the AI/ AS and how they have been verified,

(3) for accident investigators, secure storage of sensor and internal state data, comparable to a flight data ²recorder or black box.)



3.4 Multidimensional AI Models for Qualities (quantified), And Costs



Principle 4 - Transparency

Level: Organization, Status: Not Determined Type: Value, Labels: no labels Edit

Is Part Of:
GENERAL IEEE PRINCIPLES

Consists Of: i Explicability interpretability interpretability





Example Narrow Priority Objective.Wish [Decision Criteria = Decision Reasons, Transparency Level = {Data Trustworthiness, Software Engineering Process Transparency, Technical Inspection Transparency), Stakeholder Type = Accident Investigator, A/IS Incidents = {Threatening Behaviour, Operator Failure}, Incident Severity = Lethal Force Activation, A/IS Operation Types = Cyber Defence] @ 24 May 2022 : 42 <- Tsg. Just an initial narrow parameters example. 22 uantification

Ambition Level: transparent A/IS are ones in which it is possible to discover how and why a system made a particular decision, or in the case of a robot, acted the v/t 2 of v Stakeholders: A/IS TECHNOLOGISTS, CONSUMER ADVOCATES, ENGINEERS, ETHICISTS, GOVERNMENTS, IEEE P7001 Transparency Of Autonomous Systems Star Apply DUALS INTERNA. Scale: % of [A/IS Operation Types] where the [Decision Criteria] has a [Transparency Level] for a [Stakeholder Type] for [A/IS Incidents] of [Incident Sevential] Example Narrow Priority Current Level. Status: 24 [Stakeholder Type = Accident Investigator, A/IS Incidents = {Threatening Behaviour, Operator Failure}, Incident Example Narrow Priority Objective. Wish: 42 [Decision Criteria = Decision Reasons, Transparency Level = {Data Trustworthiness, Software Engineering Process Transparency, Technical Inspection Transparency}, Sta.

Relations: IEEE STANDARD P7001 - Transparency Of Autonomous SystemS, Principle 1 - Human Rights.

Issue: IssueActionIssue: How can we ensure that A/IS are transparent?

IEEE2.Rationale: Transparency is important to each stakeholder group for the following reasons:1. For users, transparency is important because it provides a simple way for them to understand what t.

Candidate Recommendation IEEE2 .Implementation Plan: Candidate RecommendationDevelop new standards* that describe measurable, testable levels of transparency, so that systems can be o...

IEEE STD.Relations: IEEE STANDARD P7001 - Transparency Of Autonomous SystemS, Principle 1 - Human Rights.

Candidate Recommendation Develop new standards that describe measurable, testable levels of transparency, so that systems can be objectively assessed and levels of compliance determined.

> **DUANTee**1 Kenner's Serverdan & Camper-How to quantify values. fors #. CB .. mes Airt makes, ENDought

Free copy. https:// tinyurl.com/ Quanteer

methods

21

22



3.5 Multidimensional AI Models for Qualities (quantified), And Costs Main definition and

Quantification: details of dimensions

Detail of the Tran

Ambition Level: transparent A/IS are ones in which it is po _____le to discover how and why a system made a particular decision, or in the case of a robot, acted the way it did.

Stakeholders: A/IS TECHNOLOGISTS, CONSUMER AD OCATES, ENGINEERS, ETHICISTS, GOVERNMENTS, IEEE P7001 Transparency Of Autonomous Systems Standard, INDIVIDUALS, INTERNATI...

Tag.Scale:

% of [A/IS Operation Types] where the [Decision Criteria] has a [Transparency Level] for a [Stakeholder Type] for [A/IS Incidents] of [Incident Severity]

Templates 👻

A/IS Incidents: defined as:

Accidents, Misuse, Confusing Communication, Inexplicable Action, Threatening Behaviour, Financial Loss, Equipment Failure,

A/IS Operation Types: defined as:

Security Operation, Recovery Operation, Weapon Operation, Road Transport Operation, Flight Operation, Flight Monitoring. Operation, Flight Monitoring. Cyber Attack, Cyber Defence,

Decision Criteria: defined as:

Right To Explanation, Decision Challenge, Decision Assessment, Decision Reasons,

Incident Severity: defined as:

Lethal Force Activation, Traceability Log Fault, Violations of Systems Embedded Community Norms, Mitigation Failure, Trust V

Stakeholder Type: defined as:

Accident Investigator, Validation Instance, Certification Instance, Private Person User, Responsible Component Provider, Response Court Process, Non Court Adjudication Process, Arbitration Process, Public Confidence,

Transparency Level: defined as:

Trusted User Level, Data Trustworthiness, Algorithmic Trustworthiness, Situational Awareness, Trace Records, Intentional Dece Consequences of Action or Decision, Ordinary Language, Source Code, Transparency Respecting Security, Software Engineering Process Transparency, Technical Inspection Transparency,

> The example shows the use of [Scale Qualifiers] to define a quality at a high level And to enable us to model very complex systems, with any combinations of [Scale Qualifiers]

levels of compliance determined.



Scale: % of [A/IS Operation Types]

where the [Decision Criteria]

has a [Transparency Level]

for a [Stakeholder Type] for [A/IS Incidents] of [Incident Severity].



The result of this is that Qualities are really well defined We can model very complex Al systems We can understand Al at a more-detailed level

Qualifiers like [Transparency Level]

- Al System Quality: Transparency, Scale
 - Al System Quality: Transparency, Scale <- identifying a Qualifiers like [Transparency Level] sub-dimension

<- Defining a Scale sub-dimension In terms of a set of Conditions

Free copy. Explains details https://tinyurl.com/Quanteer

3.7 Multidimensional AI Models for Qualities (quantified), And Costs

Principle 2 – Prioritizing Well-Being

Level: Organization, Status: Not Determined Type: Value, Labels: no labels Edit

Is Part Of:
GENERAL IEEE PRINCIPLES

Wish [Human Type = {Child. Teen, Female, Strong, Poor}, Satisfaction Type = Largely Satisfied, Life Type = Student, Conditions Of Life = {Health, Education}, Balance Effect Ratio = Positive Effects Dominate] @ 09 Apr 2029 : 42 <- Tom Gilb. Random selection of Scale Parameters to illustrate the use of them to define levels, and dimensions.

Ambition Level: encompassing human satisfaction with life and the conditions of life as well as an appropriate balance between positive and negative affect.

Stakeholders:

Scale: % [Human Type] [Satisfaction Type] with [Life Type] and [Conditions Of Life] and a [Balance Effect Ratio]

[Human Type = {Child. Teen, Female, Strong, Poor}, Satisfaction Type = Largely Satisfied, Life Type = Student, Conditions Of Life Status: 10

Wish: 42 [Human Type = {Child. Teen, Female, Strong, Poor}, Satisfaction Type = Largely Satisfied, Life Type = Student, Conditions Of Life = {Health, Education}, Balance Effect Ratio = Positive Effects Dominate.

Note: This entire arbitrary construction, of the Well-Being specification is not intended as a correct or complete, or agree definition. On the contrary, it is intended to show1. paths we mig...

Free copy. https://tinyurl.com/Quanteer **Explains these quantification methods**

>

3.8 Multidimensional Al Models for Qualities (quantified), And Costs

DETAIL OF Scale for Prioritising Well Being

Ambition Level: encompassing human satisfaction with life and the conditions of life as well as an appropriate balance between positive

Stakeholders:

a Tag.Scale:

% [Human Type] [Satisfaction Type] with [Life Type] and [Conditions Of Life] and a [Balance Effect Rat

Templates -

Balance Effect Ratio: defined as:

Positive Effects Dominate, About Equal Positive and Negative Effects, Negative Effects Dominate

Conditions Of Life: defined as:

Health, Family, Wealth, Security, Education, Mobility, Freedom of Religion, Freedom of Opinion, Fr Opportunity, ...

Human Type: defined as:

Child. Teen, Adult, Retired, Male, Female, Other Sex, Strong, Weak, Rich, Poor, Ha

Life Type: defined as:

Student, Working, Retired, Volunteer, In Care

Satisfaction Type: defined as:

Extremely Satisfied, Largely Satisfied, Dissatisfied, Extremely Dissatisfied

Target Time Units:

Calendar Date

Advanced...

V

Source: by tomailb - Apr 30th 2019 00:26

Candidate Recommendation

Develop new standards that describe measurable, testable levels of transparency, so that systems can be objectively assessed and levels of compliance determined.

##+e for	Cluanteer QUANTEER How to quantify values.
eedom of Association, Equal	
	гее сору.

Well Being' Scale of measure definition

Quanteer **Explains these** quantification methods

https://

3.10 'ACCOUNTABILITY Multidimensional AI Models for Qualities (quantified), And Costs

Ambition Level: Based on the cultural context, application, and use of A/IS, people and institutions need clarity around the manufacture and deployment of these systems to

Issue: IssueAction->How can we assure that designers, manufacturers, owners, and operators of A/IS are responsible and accountableMaybe see Candidate Recommendat

Stakeholders: CULTURE, INDIVIDUALS, ORGANIZATIONS.

a Tag.Scale:

% of cases of [Harm] where we can establish [Responsibility] towards [Peop

Templates -

A/IS Application: defined as:

Transport, Advice, Medical, Education, Household, Care Institutions, ...

Harm: defined as:

Potential Harm, Actual Harm, Severe Harm, Mass Harm, Fatal Harm, Psych

People And Institutions: defined as:

Individual Stakeholders, Institutional Stakeholders, Innocent Victims, ...

Responsibility: defined as:

Legal Culpability, Ethical Responsibility, Potential Responsibility, Acknowledged Responsibility, Indirect Responsibility, Financial Responsibility, Mitigation Responsibility

 \sim

Target Time Units:

Calendar Date

le And Institutions] for a given [A/IS Applicatio	n] Quantification: definition and dimensions
	Free for formatting, to add qualifiers
nuntahility Scale detail	https://tinyurl.com/Quanteer
	Explains these quantification met
	IF YOU CANNOT LEARN THIS,
	MOVE ASIDE, ChatGPT is
nological Harm	Very good at it. TG 2024
lere great that the	OUANTeer
	How to quantify values.
	Constant Con
	Alari az dos, EXitrandos Ballog case dote Castrado Version est Data la Heritadiag be acturary
	to the second second to part on out to part of the description of the second to be a second to b

Advanced...

Summary thus far

All Al-system <u>qualities</u>, and other stakeholder values

Here is a free detailed process for doing that quantification Without exception https://tinyurl.com/Quanteer

can be defined as - a quantity, a metric'.

So that we have a clearer, stable, tool for comparing - Al-system performance.

Quanteer

How to quantify values.

https://LeanPub.com/Quanteer

6 By Tom Gilb, 2020-2021

Tom. @ Gib. com

Start 120820, END 240820

Edit 04 June 2021

Current Version see Date in File Heading for accurate current version

mid 10th partiany, from Little scientus they good, hey maph?

'Kvantør' =. Norwegian for Quanteer.

31

4.0 The AI Stakeholder Model

- * Another Technoscope'** for understanding an Al-system, is by analyzing, and listing, all critical, or interesting, stakeholders.
- * Stakeholders are the direct sources of requirements for the Al-system
- * against which, various Al-systems can be evaluated

* Is an Al-system the best for your entire set of critical stakeholders?

** FREE COPY BOOK : TECHNOSCOPES. https://www.gilb.com/offers/YYAMFQBH/ 32

4.1 The Al Stakeholder Model

* Another Technoscope for understanding an Alsystem is by analyzing and listing all critical or interesting stakeholders.

IEEE P7000 Certification Program

* These are the direct sources of requirements for the Al-system

* Against which various Alsystems can be evaluated

4.2 The AI Stakeholder Model (a tool for understanding the AI Black Box)

IEEE AI-Ethics Standard Stakeholder Analysis (in ValPlan app)

4.3 The Al Stakeholder Model

Why does a stakeholder model help us to understand the Al-system?

- It helps you discover stakeholder needs, and specific stakeholder levels, and types (who, what, where) of those needs (not just a generic one-size-fits all needs)
- 2. It helps evaluate the Al system from the point of view of the specific stakeholders
- 3. It gives you a chance to start an early evolutionary incremental value flow, towards those stakeholders.
- 4. It helps identify 'hard constraints', like laws and budgets, which would make selection of a particular Al-system invalid or a big mistake.

5. Understanding AI Strategies in light of Multiple Quality and Cost Attributes

* We can estimate the effectiveness and costs of Al-systems against all value objectives and costs

* Because, if we quantify our Al Stakeholder Values,

* we have some clear idea of what we need to choose; in terms of an Al-system

EVALUATE AN AI TECHNOLOGY'S QUALITIES From Level: Stakeholder To Level: Product O Help mel + Add -10 O Settings Ideal XAI Technol... Notional Technique UC Berkeley Deep ... UCLA Pattern Theory+ Requirements Al Accountability Status: 1 -> Wish: 90 % of [AI S ... 56 % 2222 2222 % of [Al Scenarios] where [Al Suppliers] Al Controllability 0 Status: 50 -> Wish: 98 % of [AI S... 88 % 2222 0 2222 6 of [Al Situations] and [Al Accidents] Al Ethicality Status: 5 > Wish: 90 % #Respect... -4 % ???? ???? % Respect Violations of a Respect Viola Al Privacy Status: 1 > Tolerable: 0.5 % of [Priv... 80 % ???? ???? % of (Privacy Breach Types) annually for Al Robustness Status: 70 > Wish: 99.9 % #intende. 33 % ???? ???? % Intended Uptime for (Al Consumers) [AI Al Safety Status: 100 > Wish: 10 Number of ... 47 % ???? ???? Number of annual AI Damage occurrences Al Transparency Status: 10 > Wish: 99.998 % of selec. 47 % 22 % 2222 % of selected [AI Cases] where we can [U Al Usability 33 % Status: 1 > Wish: 0.1 Speed in M.. ???? ???? Speed in Minutes for an [Al User Type] t Al Learning Performance Status: 0 > Wish: 50 % [Learnin... 84 % ???? ???? % [Learning Performance] improvement fro 22 % 0 % 464 % 900 % 0 36 im Of Values: Operational Cost Budget: 0 > Status: 100 % of a def. 89 % 99 % 99 % ???? % of a defined Budget m Of Development Resources: 99 % 0 96 99 % 99 96 0 % Value To Cost: 0.20 4.79 minite Ratio (Worst Case) 0.10 0.00 Ratio (Cred. - adjusted) Ratio (Worst Case Cred. - adjusted)

FREE: Value Impact Estimation BOOKLET, https://tinyurl.com/VIEbooklet VIE BOOK PDF FOLDER, 2022, 70 PAGES. First draft complete

* We can estimate the effectiveness and costs of Al-systems against all value objectives and costs

* Because, if we quantify our Al Stakeholder Values,

* we have some clear idea of what we need to choose; in terms of an Al-system



6. Teaching The Multi-AI model to students and research students as a tools for researching AI developments. AN 'ACADEMIC' VIEWPOINT

* Teach 'Planguage' as a basic tool for 'life', and all Al research.

* See https://www.researchgate.net/publication/237129623_Competitive_Engineering_A_Handbook_for_Systems_Engineering_A_Handbook_for_Systems_Engineering_ring

* and get free copy of 'Competitive Engineering'

* Use Planguage to analyze Al systems and write papers and thesis

* To Be Clear: I am not here to sell books or services, I am retired, and I want to convince you to try my ideas, and I want to give you free written materials to teach them.





7. The Quality distinctions between Large Language Model (LLM) AI and the next AI generation (AGI, Artificial **General Intelligence**)

- and fake news.
- * Artificial General Intelligence: is Not based on Big Data.
 - AGI
 - papers. (Disclaimer: I am investor here)
 - * There are some videos on site to explain
 - * Peter Lawrence explains this (AGI technicalities):

https://www.mygreatlearning.com/

Artificial Narrow Intelligence (ANI)



Stage-1

Machine Learning

and solves one proble



Artificial General Intelligence (AGI)



Stage-2

Machine Intelligence

that is as smart as a

Artificial Super Intelligence (ASI)



Stage-3

Machine Consciousness

practically every field

* Large Language Models: merely process data and spit out tailored versions of it, like 'magic'. Often wrong

* They are based on teaching Al to 'reason like a human' and solve interesting problems. The Al future is

* See graphmetrix.com for emergence of a real product that already does this on building drawings and

* https://medium.com/@peter.lawrence_47665/knowledge-graphs-large-language-models-theability-for-users-to-ask-their-own-questions-e4afc348fa72



AN AI-PRODUCER OPINION

My one word response to this ("The AI Dilemma") is "delusional".

The fact is, there is no Intelligence at all in LLMs,

and because of this, they can't do the intelligent things people think they can

because they use words that were output from intelligence.

There is A for sure.

The primary danger of these things is fooling people, so a **big productivity boost for criminali**ty.

Outside of that, LLMs are literally pretty useless for anything mission critical that needs to work autonomously.

So these guys are crazy in thinking that anyone will lose their jobs, the whole economy will change etc. For sure the criminals will be able to commit more crime more easily however.

Those radiologists that were supposed to be out of work by now... those truck drivers, taxi drivers... etc etc. Its crazy hype where billions will be lost as has happened with autonomous driving.

Intelligence is needed for these to work, and that is fully missing.

(Except for what we are working on ... ;) <- GRAPHMETRIX.COM

Fred Gibson Founder & CEO mobile: 415.335.8232

APRIL 11 2023

https://vimeo.com/809258916/92b420d98a The Al Dilemma

that le ck



https://rocketreach.co/frederick-gibsonemail_7244557



AN AI-PRODUCER OPINION

Email 14 April 2023

I would say that the most important point is that LLMs and ml/dl have no intelligence at all,

and naming them "Artificial Intelligence" is both <mark>false advertising</mark> and promising something the systems can never deliver, intelligence.

They can get you to mission-critical (and very useful) perception *if* the input is pointed at the world of space-time,

that is measurements that an intelligent system can then use to reason with,

but pointing the input at words bypasses this perception potential entirely and gives you just word pattern outputs of random usefulness.

Fred Gibson Founder & CEO mobile: 415.335.8232

ml=Machine Learning, dl= Deep Learning, LLM= Large Language Model

The Al Dilemma



How do we measure to see it Artificial Intelligence is acting like a human?

Even if we reach that state where an Al can behave as a human does, how can we be sure it can continue to behave that way?

We can base the human-likeness of an Al entity on the:

- Turing Test
- The Cognitive Modelling Approach
- The Law of Thought Approach
- The Rational Agent Approach

TEST THE WITH THESE TESTS



8. The Penta Model as a basic high level view of any Al system. * The Penta Model is a high-level simplification of the analysis ideas in Planguage (Competitive Engineering, gilb.com), https://tinyurl.com/PentaPaper

Funct Tax: Compute all relevant Taxes

Design Constraint A: Avoid Patented Designs.

Function Constraint P: Avoid Political Bias

Strategy 1: No Cure No Pay Contract

Design D: Digital Twin

Designs

Budget: € 1 million

Deadline: : January 1 Next Year



9. The Principles of Al Understanding, outside the blackbox.

- and costs

- gives opportunities for early results delivery, to selected segments.
- more knowledge, by experimentation, testing, and pilot use.

© tom@Gilb.com 2023. Version 18 April 2023

1. We cannot hope to translate AI detailed complexity into understanding of qualities

2. We have to work outside the Al-Black-Box with multiple quantified qualities and costs.

3. We can compare Al-systems, against their known ability, to satisfy quality objectives and cost constraints, by using an Impact Estimation Table. https://tinyurl.com/ VIEbooklet. And we can use AI-LLMs to draft stakeholders, values, & costs, tables.

4. Identifying critical Al stakeholders gives better understanding of Al requirements, and

5. If you cannot get real evidence about Al-system qualities and costs, then you have got specific 'known unknowns', risks. Which (the known unknowns) might be made into





- * tom@gilb.com
- * Www.gilb.com
- * <a>www.graphmetrix.com (Investor, Advisory Board)
- * @ImTomGilb
- * http://www.linkedin.com/in/tomgilb
- * twitter.com/imtomgilb
- * Location: Near Oslo Norway
- * Retired (83), but enjoy discussing and sharing
- * Honorary Fellow of BCS (2013)
- https://www.researchgate.net/profile/Tom-Gilb
- * https://leanpub.com/u/tomgilb

LAST SUDE



SLIDES FOLDER



I published 'Software Metrics' in 1976 - and been fighting for a metrics (engineering) culture, ever since.

Al, Security, and IT need metrics badly.

I am only 83 (July 2024), so I will irritate 'poetic methods' people for guite a while

And after that my writings will fight on!





XAI Talk 2019 Slides

You probably do not have time in a 45 minute lecture to go deeper BUT

Here is more technical detail on quantifying Al qualities

https://www.researchgate.net/publication/ 381932584 eXplaining Al A Serious 'Multidimensional Metrics Attack on Poor Al 'Academic and Stand ards' Thinking Planning





'XAI' "eXplaining A.I." A Serious 'Multi-dimensional Metrics Attack' on **Poor AI 'Academic and Standards' Thinking & Planning**





'Principles' = 'Qualities' ! With the usual lack of quantified definition

Oct. 2016 : "Preparing for the future of A!" from USA





© 2017 International Business Machines Corporation

Proposal of Discussion toward Formulation of AI R&D Guideline

Referring OECD guidelines governing privacy, security, and so on, it is necessary to begin discussions and considerations toward formulating an international guideline consisting of principles governing R&D of AI to be networked ("AI R&D Guideline") as framework taken into account of in R&D of AI to be networked.

Proposed Principles in "AI R&D Guideline"

- 1. Principle of Transparency
- 2. Principle of User Assistance

make choices

3. Principle of Controllability

Ensuring controllability of the Al network system by humans

- 4. Principle of Security
- 5. Principle of Safety
- 6. Principle of Privacy
- 7. Principle of Ethics
- 8. Principle of Accountability

https://obamawhitehouse.archives.gov/sites/default/files/whitehouse_files/microsites/ostp/NSTC/preparing_for_the_future_of_ai.pdf

Ensuring the abilities to explain and verify the behaviors of the AI network system

Giving consideration so that the AI network system can assist users and appropriately provide users with opportunities to

Ensuring the robustness and dependability of the Al network system

Giving consideration so that the AI network system will not cause danger to the lives/bodies of users and third parties

Giving consideration so that the AI network system will not infringe the privacy of users and third parties

Respecting human dignity and individuals' autonomy in conducting research and development of AI to be networked

Accomplishing accountability to related stakeholders such as users by researchers/developers of AI to be networked

IBM

Al Accountability

Steps to Accountability

©Eve Ash eveash.com

https://www.smartcompany.com.au/people-human-resources/the-a-listsix-steps-to-accountability-in-your-business/



Al Accountability Defined Fuzziness -> Quantified. Structured, Enriched



Al Accountability Scale **Detailed Scale Parameters** Why is this Useful?



Al Controllability

what are other words for controllability?



manageability, handleability, docility, tractability, governability, maneuverability, directional control



Thesaurus.plus

Al Controllability: Quantified



The current level of 'Controllability' Quantified

Rich Slice Example.Status [Al Situations = Potential Attack Detected,Attack in Progress, Al Accidents = All, AI Consumers = Medical, Government, AI Scenarios = Life Critical Decisions, Health Decisions, AI Management = All, Al Status = Normal Successful Operation] @ 09 Jan 2019 : 50 ±49 ?? SWAG <- tsg example

Scale: % of [Al Situations] and [Al Accidents] for [Al Consumers] and [Al Scenarios] and [Al Situations] where Successful Al Control can be applied by [Al Management] resulting in [Al St. Rich Slice Example. Status: Level: 50 ±49 ?? SWAG [Al Situations = Potential Attack Detected, Attack in Progress, Al Accidents = All, Al Consumers = Medical, Government, Al Scenarios = Life Cr.

AICK to show the Sloebar (press alt+s

>

The desired and valued level of 'Controllability'

Al Controllability: Scale of Measure Detail

Ambition Level: "Ensuring the controllability of the AI network system by humans"

Stakeholders: AI SYSTEMS, AI System Management, AI Users.

Tag.Scale:

Scale Description: 😧

Status]

Al Accidents: defined as:

Al Hardware Failure, Al Software failure, Al Data Failure, Al Data Result Failure, Al Action Failure, [Acci

Al Consumers: defined as:

Citizens, Medical, Transport, Government, Telecoms, Finance, Industry, Commerce, Distr Senior Homes, Private Homes, Apartments, Offices...

Al Management: defined as:

Al Supplier Service, User Organization Service, User Service, Service Organization

Al Scenarios: defined as:

Al Situations: defined as:

tual Corruption Happened, System Recovered Without Damage, No Damage Occurred

Al Status: defined as:

Normal Successful Operation, Dominant Learning Mode, Partial Decision Failures, Frequent Decision Failures, Out of Operation, Recovery Mode, Secondary Backup System, Mode, Substantial Human Intervention Necessary



Al Ethicality



https://medium.com/salesforce-ux/how-to-build-ethics-into-ai-part-i-bf35494cce9



AI Ethicality: Summary and Ethicality Scale of Measure

6 Dimensions

St Control Con	h [Respect Violation Type = Legal Violation rernment, Industry, Commerce, Autonomo isions, Al Suppliers = All, Market Area = C real		
Ambition Level: "Respecting human digni	and individuals' autonomy in conducting r		
Stakeholders: Al Human Tuners, Al Syster	Management, Al Users, Al VICTIMS.		
Tag Scale:			

Tag.Scale:

Scale Description: 🕜

% #Respect Violations# of a [Respect Violation Type] towards [AI Consumers] in [

Al Consumers: defined as:

Citizens, Medical, Transport, Government, Telecoms, Finance, Industry, Commer Senior Homes, Private Homes, Apartments, Offices...

Al Scenarios: defined as:

Freeware, Paid Services, [Life Critical Decisions], Health Decisions, Large Sca

Al Suppliers: defined as:

App Developers, Data suppliers, Academic Research, Startups, Consultants, IO

Market Area: defined as:

Core Initial Market, Worldwide, Europe, Asia, China, India, Africa, USA, Latin Am

Respect Violation Type: defined as:

Legal Violation, Policy Violation, Guideline Violation, Rule Violation, Standards violation, Economic Violation, Safety Violation, Privacy Violation, ...

Respect Violations: defined as:

Any Al action or consequence which is in fact perceived by individuals or groups as negative, and which could somehow be avoided by system change

			Ċ
on,Star	ndards Violat nicles: Al Sce	tion,N	legative PR Violation, Al Consumers = s = Health Decisions.Entity Survival
ore In	itial Market]	@ 09	Jan 2027 : 90 ±9)) <- tsg example only
esearc	h and develo	pmen	t of AI to be networked."
			by tomgilb - Jan 9th 2019, 12:28 🗪 0 🖆 🛨 🛅 🏥
Al Sc	ي ک		ppliers] Annually in a given [Market Area]
ce, D	DEFI cale P	Str	Entertainment, Autonomous Vehicles, IOT Consumers, Military, Police,
le De	NITIO	ucture	cisions, Entity Survival Decisions, Political Choices, Games
sup	eters		tions,
erica,			
			Jative PR Violation, Gender Violation, Racial Violation, Religious Vi-
	nativo pad w	hich -	rould comphaw be avaided by system shance

Al Ethicality: Scale Parameter Detail

Tag.Scale:

Scale Description: 😮

in a given [Market Area]

Al Consumers: defined as:

Citizens, Medical, Transport, Government, Telecoms, Finance tonomous Vehicles, IOT Consumers, Military, Police, Senior

Al Scenarios: defined as:

Freeware, Paid Services, [Life Critical Decisions], He Decisions, Political Choices, Games...

Al Suppliers: defined as:

App Developers, Data suppliers, Academic Research, Start

Market Area: defined as:

Core Initial Market, Worldwide, Europe, Asia, China, India, Afri

Respect Violation Type: defined as:

Gender Violation, Racial Violation, Religious Violation, Economic Violation, Safety Violation, Privacy Violation, ...

Respect Violations: defined as:

Any Al action or consequence which is in fact perceived by individuals or groups as negative, and which could somehow be avoided by system change



Al Ethicality: Summary of the 'WISH LEVEL' Status Wish 90 5

Wish [Respect Violation Type = Legal Violation, Standards Violation, Negative PR Violation, Al Consumers = Government, Industry, Commerce, Autonomous Vehicles, Al Scenarios = Health Decisions, Entity Survival Decisions, Al Suppliers = All, Market Area = Core Initial Market] @ 09 Jan 2027 : 90 ±9)) <- tsg example only not real

Here is an example of selecting a *prioritized*, <u>critical</u>, part of the action, maybe in practice 1/50 of total budget, and 1/50 total time to deadline (a 'sprint' if you like) Value delivered early, design experiments possible, feedback and correction possible quickly FREE GIFT REVIEW COPY FOR YOU ALONE. NO COUPON CODE **REQUIRED.**



>

Decompose



Al Learning Performance

Al Learning Performance: Overview of a structured quantified requirement



stakeholder value



Al Learning Performance: Wish Level



Al Learning Performance: **Status Level**

Stakeholders: ACADEMIC INSTITUTIONS, AI Human Tuners, INDUSTRIAL AI PARTNERS. Show Sideba by tomgilb - Jan 10th 2019, 03:39 👳 0 % [Learning Performance] improvement from a defined [AI Technique Performance Baseline] for [AI Cases] and [AI Scenarios] and for [AI Consumers] Time Units: Calendar Date Al Cases: defined as: Automated Decision, Prediction for Humans, Decision for Machine System Action Al Consumers: defined as: Citizens, Medical, Transport, Government, Telecoms, Finance, Industry, Commerce, Distribution, Academic, Entertainment, Autonomous Vehicles, IOT Consumers, Military, Police, Senior Homes, Private Homes, Apartments, Offices... Al Scenarios: defined as: Freeware, Paid Services, [Life Critical Decisions], Health Decisions, Large Scale Decisions, Long Term Decisions, Entity Survival Decisions, Political Choices, Games... Al Technique Performance Baseline: defined as: Estimated Learning Performance, Measured Learning Performance, Unknown Learning Performance, ... Learning Performance: defined as: Perfect Recognition, Perfect Decision, Extremely Safe Levels, Extremely Fast Lean. Domain, Extremely Fast Learning New Real Time Information, ... Source: Defining a Scale of tsg Measure for Status: 'Learning Performance' 0 baseline 0 [Learning Performance = Extremely Safe Levels, Al Technique Performance Baseline = Estimated L by means of [Scale Entity Survival Decisions, AI Consumers = Autonomous Vehicles] (to be delivered by end of: 10 Jan 2019) Parameter] conditions

Scale:

Al Learning Performance: Darpa Source Level Spec: example of source.







/ERS TIZENSH DER 10 **1**12/03 GEI BAN 411.40 AGE SCHOO INDER





<- tsg swag

Stakeholders: ACADEMIC INSTITUTIONS, AI DATA SOURCES, AI Users, AI VICTIN

S	ale Description: 😧
ſ	% of [Privacy Breach Types] annually for [Al Consumers] in defined [A
	Al Consumers: defined as:
	Citizens, Medical, Transport, Government, Telecoms, Finance, Industry, Con Senior Homes, Private Homes, Apartments, Offices
	Al Situations: defined as:
	[Situations], Potential Attack Detected, Attack in Progress, Real Attack Detected tual Corruption Happened, System Recovered Without Damage, No Damage Occ
	Market Area: defined as:
	Core Initial Market, Worldwide, Europe, Asia, China, India, Africa, USA, Latin Ame
	Privacy Breach Types: defined as:
	Identity Publicized, Financials Hacked, Personal Data Corrupted, Unauthorized A

Al Privacy



Al Safety According To Tesla, The Model 3 Is The Safest Car On The Planet

OCT 9, 2018 / BY JAY TRAUGOTT IN TESTING 🤊

And the number two and three spots go to ...

Last month it was announced the new Tesla Model 3 earned a perfect five-star crash test rating in every category from the National Highway Transportation Safety Administration. Most impressive. And now Tesla says its new mid-size sedan has the lowest probability of causing injury for any vehicle ever tested by the US government agency.

"Not only has Model 3 achieved a perfect 5-star safety rating in every category and sub-category, but NHTSA's tests also show that it has the lowest probability of injury of all cars the safety agency has ever tested," Tesla wrote in a blog post. This same post also explains why the Model 3 is so safe.



Al Safety Is your family safety worth quantifying? or is 'Make My Family Safe Again' Good enough ?



	Show Sidebar
Situations = Real Attack Detected, Actual Corruption 10 <- tsg draft 1	Happened,
to the lives/bodies of users and third parties"	
ations] in [Market Area].	by tomgilb - Jan 9th 2019, 14:32 🔍 0
Al Action Failure, [Accident]	
onal or family relationships, arguably caused directly or i	ndirectly by an Al system, even the Al is op-
ed, Potential Data Corruption, <u>Potential</u> Logic <u>Corruptic</u> rred	on, Potential Hardware Corruption, Actual
merica,	

Al Safety So here is an example of quantification of Safety. Why don't the AI Standards groups INSIST on quantification of Critical AI Attributes? Maybe... they do not intellectually understand 'quantification'?

AI Safety	
Level: Stakeholder, Type: Value, Labels: - Edit Is Part Of: DCRITICAL AI QUALITIES	
Status Wish 100 10 Mich Kill Andian Entire Al Offention Dedational Detected At a Detected At a diameter of the second state	
Wish [Al Accidents = Al Action Failure, Al Situations = Real Attack Detected, Actual Corruption Happened, Market Area = Worldwide] @ 09 Jan 2022 : 10 <- tsg draft 1	
Ambition Level: "Giving consideration so that the AI network will not cause danger to the lives/bodies of users and third parties"	i
Stakeholders: Al Human Tuners, Al Users, Al VICTIMS, Unaware People.	=
Scale: Number of annual Al Damage occurrences a result of [Al Accidents] in [Al Situations] in [Market Area].	
Status: 100 huge, no data, guess [Al Accidents = Al Action Failure, Al Situations = Real Attack Detected, Actual Corruption Happened, Market Area = Worldwide] (to be delivered by end of: 09 Jan 2019) Source: tsg draft 1	=
Wish: 10 [Al Accidents = Al Action Failure, Al Situations = Real Attack Detected, Actual Corruption Happened, Market Area = Worldwide] (to be delivered by end of: 09 Jan 2022) Source: tsg draft 1	=

Al Transparency

 Can entities train a model without having to disclose their data?

The Privacy Problem

The Transparency Problem

 Can the activity and behavior of an AI model be transparently available to all parties without the need of trusting a centralized authority?


Al Transparency

evel: Stakeholder Type: Va		
Is Part Of: Diakenolder, Type. Va	UALITIES	
	Status	Wish
		99.998
	Status [Al Cases = Dec signed,Certify As Safe, sumers = Autonomous	ision for Machine S Al Accidents = Al S Vehicles, Market Are
	Status [Al Cases = Dec signed,Certify As Safe, sumers = Autonomous example, no data <- tsg	ision for Machine S Al Accidents = Al So Vehicles, Market Are g draft 1 example
Ambition Level: Ensuring the a	Status [Al Cases = Dec signed,Certify As Safe, sumers = Autonomous example, no data <- tsg	ision for Machine S Al Accidents = Al So Vehicles, Market Are g draft 1 example eb of the Al net
Ambition Level: Ensuring the a Stakeholders: ACADEMIC INS	Status [Al Cases = Dec signed,Certify As Safe, sumers = Autonomous example, no data <- tsg abilities to explain and verify the be	ision for Machine S Al Accidents = Al So Vehicles, Market Are g draft 1 example of the Al net System Managemer
Ambition Level: Ensuring the a Stakeholders: ACADEMIC INS Scale: % of selected [AI Cases	Status [Al Cases = Dec signed,Certify As Safe, sumers = Autonomous example, no data <- tsg bilities to explain and verify the be TITUTIONS, Al Human Tur 2, Al	ision for Machine S Al Accidents = Al So Vehicles, Market Are g draft 1 example of the Al net System Managemen ecisions] for [Al Acc

Tags on various scale levels are useful, since we can have a many levels of each type (Wish etc=) as we need

1

 Action
 Land Al Decisions = Confirm Logic as De

 ftwar
 In Damage = Potential or actual damage, Al Con

 Initial Market,Europe,USA] @ 09 Jan 2019 : 10 ±10 ??

work system
t, AI Users, Research Institutions.
dents] and consequent [AI Damage] with specific [AI Consumers] in a [Market Area].
ction, Understand AI Decisions = Confirm Logic as Designed,Certify As Safe, AI Accidents = AI Software fail...
for Machine System Action, Understand AI Decisions = Confirm Logic as Designed,Certify As Safe, AI Acci.

AI Transparency

Stakeholders: ACADEMIC INSTITUTIONS, AI Human Tuners, AI System Management, AI Users, Research Institutions.

Scale:

[Al Consumers] in a [Market Area].

Time Units: Calendar Date

Al Accidents: defined as:

Al Hardware Failure, Al Software failure, Al Data Failure, Al Data Result Failure, Al Actic

Al Cases: defined as:

Automated Decision, Prediction for Humans, Decision for Machine System Action

Al Consumers: defined as:

tonomous Vehicles, IOT Consumers, Military, Police, Senior Homes, Private Homes, Apartments, Offices...

Al Damage: defined as:

or indirectly by an Al system, even the Al is operating as designed and intended.

Market Area: defined as:

Core Initial Market, Worldwide, Europe, Asia, China, India, Africa, USA, Latin America,...

Understand AI Decisions: defined as:

Path for Analysis, ...

Status: Level: 10 ±10 ?? example, no data [Al Cases = Decision for Machine System Action, Understand Al Decisions = Confirm Lo + - : iii



Al Usability



Al Usability

Stakeholders:	
Stakeholder 🕈	Roles
From: A Users	
From: <u>Aware People</u>	
Scale: Speed in Minutes for an [Al User Type] to Se	uccessfully Use AI for [AI Scenarios] in a [Market Area].
Time Units: Calendar Date	
Al Scenarios: defined as:	
Freeware, Paid Services, [Life Critical D	ecisions], Health Decisions, Large Scale Decisions, Long T
Al User Type: defined as:	
Unaware of AI, Conscious of AI,	
Beginner, Self Taught, Safe User, Expert I	Level
Market Area: defined as:	
Core Initial Market, Worldwide, Europe, A	sia, China, India, Africa, USA, Latin America,
Successfully Use Al: defined as:	
Use AI as intended, without delays or fau	ilts.
Source:	
tsg draft 1	
Status:	
1 ?? [Al User Type = Expert Level, Al Scenarios	= Life Critical Decisions, Entity Survival Decisions, Market Area =

Source:

rs and appropriately provide users with opportunities to make choices." by tomgilb - Jan 9th 2019, 01:45 • 0 Notes Show Sidebar by tomgilb - Jan 9th 2019, 15:38 • 0

ale Decisions, Long Term Decisions, Entity Survival Decisions, Political Choices, Games...





Al Usability

Show S	idebai	T
narios = Life Critical Decisions, Entity Survival Decisions, Market		
and appropriately provide users with opportunities to make choices."		
by tomgilb - Jan 9th 2019, 01:45 Notes	♥ 0	8
s] in a [Market Area].	Û	
by tomgilb - Jan 9th 2019, 15:40 ecisions, Market Area = USA] (to be delivered by end of: 09 Jan 2019)	9 0	
tity Survival Decisions, Market Area = USA] When 09 Jan 2021	C	

Other Viewpoints

2020.

12



Source: https://www.businesswire.com/news/home/20170509005837/en/ServiceNow- Launches-Intelligent-Automation-Engine%E2%84%A2

5 3

Diagram: AI Stakeholders & AI 'Qualities' (Principles)



Al Hardware Al Human Tuners Al Logio Al Users Al System Management INDUSTRIAL AL PARTNERS

Diagram: Stakeholders, Qualities, Some 'XAI Techniques'









Approved for public release: distribution unlimited.

Decision Diagrams

Interactive Visualization

Bayesian Teaching



		🗎 www.darpa.mil/attachments/XAIP	rogramUpdate.pdf 🔿	
DAR	XAI Deve	Photo - Google Photos elopers (TA1)	https://www.darpa.mil/attachments/XAIProgramUpdate	ELUBENC
СР	Performer	Explainable Model	Explanation Interface	
	UC Berkeley	Deep Learning	Reflexive and Rational	
Both	Charles River	Causal Modeling	Narrative Generation	
	UCLA	Pattern Theory+	3-level Explanation	
my	Oregon State	Adaptive Programs	Acceptance Testing	
ouo	PARC	Cognitive Modeling	Interactive Training	
Aut	CMU	Explainable RL (XRL)	XRL Interaction	
	SRI International	Deep Learning	Show and Tell Explanation	
S	Raytheon BBN	Deep Learning	Argumentation and Pedagogy	
alyt	UT Dallas	Probabilistic Logic	Decision Diagrams	
An	Texas A&M	Mimic Learning	Interactive Visualization	
	Rutgers	Model Induction	Bayesian Teaching	
		Assured for mublic veloces distributi	and the line it and	



Approved for public release: distribution unlimited.

Explainability of Learning Techniques



Relationship Diagram AI Transparency and its stakeholders



Can we evaluate Al Technology Qualities against 9 AI Qualities and some Costs?

AI Qualities Demo 2019 / Value Decision Tables / EVALUATE AN AI TECHNOLOGY'S QUALITIES

EVALUATE AN AI TECHNOLOGY'S QUALITIES

Levels Chaloshalden, Tall avel, Duaduat

Settings	+ Add -	♦ Sort ▼	Duplicate	ງ Undo	=: A	BSOLUTE Show	Help
ne!				8 <u></u>			
Requirements			Ideal XAI Technol		Y Notional Technique		¥
()→ AI Accou Status: 1 → V	u ntability Vish: 90 %	≕: of [AISΔ%;		100%		56%	
()→ AI Contr Status: 50 →	ollability Wish: 98 9	=: % of [AI S. <u>&</u> %:		100%		88%	
()→ AI Ethica Status: 5 → V	ality Vish: 90 %	=: #Respect <u>∆%</u> :		100%	-4%		
()→ Al Privac Status: 1 → T	≫y blerable: 0.	=: .5 % of (Prover:		100%		80%	
()→ Al Robus Status: 70 →	s tness Wish: 99.9	=:)%#Inten⊴i%;		100%	33%		
()→ AI Safety Status: 100 →	/ Wish: 10	≕: Number o <u>4</u> %;		100%	47	56	
()→ AI Trans Status: 10 →	p arency Wish: 99.9	=: 998 % of 9946¢.		100%	47	%.	
()→ Al Usabi Status: 1 → V	llty Vish: 0.1 S	=: peed in M <u>a</u> %:		100%	33%		
()→ Al Learn Status: 0 → V	ing Perfo Vish: 50 %	rmance =: [Learnin∆%:		100%		84%	
Sum Of Values:		Σ%:	900 %		464 %		0
-→) Operatio Budget: 0 → S	onal Cost Status: 100	=: 0% of a d <u>a%</u> :		99%		99%	
Sum Of Development Resources:2%:			99 % 99 %			0	
Value To Cos	t:		9.10		4.70		0.
					Gun	004 0	

AI Technique Evaluation using DATA with objective evidence, sources, AND uncertainty

Requirements	- Ideal XAI Technol	- Not
()→ Al Accountability ∆: Status: 1 → Wish: 90 % of [AI S ∆%: % of [Al Scenarios] where [Al Supplie [Al Suppliers = All,] 1 09 Jan 2040	89 ± 0 100 ± 0 % 0 % (x 0.0) 100%	50 ± 10 56 ± 11 28 % (
 Al Controllability ∆: Status: 50 → Wish: 98 % of [Al S ∆%: % of [Al Situations] and [Al Accident [Al Situations = Potential Att] 109 Jan 2029 	48 ± 0 100 ± 0 % 0 % (x 0.0) 100%	42 ± 0 88 ± 0 0 % (x
 Al Ethicality ∆: Status: 5 → Wish: 90 % #Respect∆%: % Respect Violations of a [Respect Vi [Respect Violation Type = Lega] 109 Jan 2027 	85 ± 0 100 ± 0 % 0 % (x 0.0) 100%	-3 ± 0 -4 ± 0 9 0 % (x -4
()→ Al Privacy Status: 1 → Tolerable: 0.5 % of [Priv _A %: % of [Privacy Breach Types] annually [Privacy Breach Types = Financ]	-0.5 ± 0 100 ± 0 % 0 % (x 0.0) 100%	-0.4 ± 0 80 ± 0 0 % (x
 Al Robustness ∆: Status: 70 → Wish: 99.9 % #Intended % Intended Uptime for [Al Consumers] [Al Consumers = Transport,Auto] 109 Jan 2025 	29.9 ± 0 100 ± 0 % 0 % (x 0.0) 100%	10 ± 0 33 ± 0 0 % (x 33%

€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€
€

ional Technique	₩ UC B						Hide Sidebar
	???? ± 0	-Select	Impact Ta	arget			
% (x 0.5) 56%	0 ± 0 % 0 % (x 0 ????	B Tag .Va Estimate:	alue Impac	:t:			
	???? ± 0	Δ	50	٢	±	10	0
%	0±0%	Actual:					
0.0) 88%	0% (x0	Δ	scale	/alu(3	±	0	0
		Credibility	y:			. Se	
	<u>????</u> ± 0		_0				
%	0 ± 0 %	Evidence					
(0.0) %	0% (x0 ????	MEASU	JRED IN PI	HD STUDY, :	3 SITES		
	???? ± 0	Source:					
%	0 ±0%	UCLA	JRL				
0.0)	0% (x0		10 de 974 de			8.040 10418 BLF	
80%	<u> </u>	tomgilb minutes	added a cor ago 💼	nment - 4 mi	nutes ago -	edited by to	mgilb - 4
	???? ± 0	тнія	SIS JUST A	STRUCTU	AL EXAMPI	E WITH	
%	0 ± 0 %	FAK	E ASSERTI	ONS TO SH	IOW THE S	TRUCTURI	E 🕜
0.0)	0% (x0	OF A	AN ESTIMA	TE. TSG			
				Pres	s #+e to sl	now editing t	toolbar.
		Ac	ld Comme	nt			

Defined Concepts: Not least for [Scale Parameters]

AI Accidents	Al Hardware Failure, Al Software failure, Al Data Failure, Al Data Result Failure, Al Action Failure, [Accident]	@tomgilb	Jan 9th 2019, 10:00	tomgilb	a
AI Consumers	Citizens, Medical, Transport, Government, Telecoms, Finance, Industry, Commerce, Distribution, Academic, Entertainment, Autonomous Vehicles, IOT Consumers, Military, Police, Senior Homes, Private Homes, Apartments, Offices	@tomgilb	Jan 9th 2019, 10:06	tomgilb	
Al Management	Al Supplier Service, User Organization Service, User Service, Service Organization	@tomgilb	Jan 9th 2019, 09:59	tomgilb	
Al Scenarios	Freeware, Paid Services, [Life Critical Decisions], Health Decisions, Large Scale Decisions, Long Term Decisions, Entity Survival Decisions, Political Choices,	@tomgilb	Jan 9th 2019, 06:09	tomgilb	Û
AI Situations	[Situations], Potential Attack Detected, Attack in Progress, Real Attack Detected, Potential Data Corruption, [Potential] Logic [Corruption], Potential Hardware Corruption, Actual Corruption Happened, System Recovered Without Damage, No Damage Occurred	@tomgilb	Jan 9th 2019, 10:00	tomgilb	
Al Status	Normal Successful Operation, Dominant Learning Mode, Partial Decision Failures, Frequent Decision Failures, Out of Operation, Recovery Mode, Secondary Backup System Mode, Substantial Human Intervention Necessary	@tomgilb	Jan 9th 2019, 10:00	tomgilb	
Al Suppliers	App Developers, Data suppliers, Academic Research, Startups, Consultants, IOT suppliers, Teaching Institutions,	@tomgilb	Jan 9th 2019, 06:08	tomgilb	t
Accountability Strength	Legally Compliant, Contractually Compliant, Standards Compliant, International Guidelines Compliant,	@tomgilb	Jan 9th 2019, 06:08	tomgilb	
ife Critical Decisions	At least one human life can be lost or effectively destroyed by bad decisions as a result of the AI tool	@tomgilb	Jan 9th 2019,	tomgilb	E

Quantitative Evaluation of AI Technology Qualities using 'Impact Estimation Tables'

EVALUATE AN AI TECHNOLOGY'S QUALITIES

From Level: Stakeholder To Level: Product

	Settings	+ Add -	🔹 Sort 🔫	Duplicate	່ງ Undo	. =: /
	lequirement	S		-₩- Ideal XAI	Technol	* Notio
	() → AI Acc Status: 1 →	Wish: 90 %	=: of [Al S∆%;		100%	
	()→ AI Co Status: 50	ntrollability → Wish: 98 9	=: % of [AI S <u>∆</u> %:		100%	10-
6	() → AI Eth Status: 5 →	icality ▶ Wish: 90 %	=: #Respec <u>h</u> %:		100%	-4%
>	() → Al Priv Status: 1 →	vacy ▶ Tolerable: 0.	=: .5 % of [P <u>r</u> iv ₆₇		100%	
	() Al Ro Status: 70	bustness ➔ Wish: 99.9	=:) % #Inte <u>nde</u>		100%	33%
	() → AI Sat Status: 100	iety) → Wish: 10	≕ Number @f%:		100%	47%
	()→ Al Tra Status: 10	nsparency → Wish: 99.9	=: 998 % of <u>aele</u> c		100%	47°
	() → AI Us: Status: 1 →	ability ♦ Wish: 0.1 S	≓: peed in M <u>a%</u> :		100%	33%
	() → Al Lea Status: 0 →	wish: 50 %	rmance =: [Learnin.A%:		100%	

	Help me!		9
nal Technique	- UC Berkeley Deep	- UCLA Pattern Theory+	* Charles Show Sideba
56%	?77?	????	????
88%	????	????	????
80%	This is a quid to show the ar AI V	e potential for eva Techniques nd their effects or alues (or 'Qualitie	ut real data: aluating Al n es')
16	3335	2555	????
%	????	????	????
	????	?? ??	7777
84%	7777	?? ??	2777

88

A more detailed example of rating 'uncertainty' (± 10) and Credibility $56\% \times 0.5 = 28\%$

Mock example

EVALUATE AN AI TECHNOLOGY'S QUALITIES From Level: Stakeholder To Level: Product A: INCREMENTAL Show Sidebar Duplicate... O Undo.. + Sort -Settings... + Add me! 🛉 Ideal XAI Technol... Yotional Technique Requirements Al Accountability 50 ± 10 89 + 0Status: 1 > Wish: 90 % of [AI S ... $100 \pm 0\%$ 56 ± 11 % Δ%: % of [Al Scenarios] where [Al Supplie.. 0 % (x 0.0) 28 % (x 0.5) [Al Suppliers = All, ...] 56% 🎁 09 Jan 2040 Al Controllability 48 ± 0 42 ± 0 Δ: Status: 50 → Wish: 98 % of [AI S ... ∆%: 100 ± 0 % 88 ± 0 % % of [AI Situations] and [AI Accident... 0 % (x 0.0) 0 % (x 0.0) [Al Situations = Potential Att...] 🛗 09 Jan 2029 ()→ AI Ethicality 85 ± 0 -3±0 Δ: Status: 5 → Wish: 90 % #Respect... Δ%: -4 ± 0 % $100 \pm 0\%$ % Respect Violations of a [Respect Vi... 0 % (x 0.0) 0 % (x 0.0) [Respect Violation Type = Lega...] -4% 🛗 09 Jan 2027 ()-) Al Privacy $\textbf{-0.5}\pm0$ Δ: -0.4 ± 0 Status: 1 → Tolerable: 0.5 % of [Priv.∆%: $100 \pm 0\%$ 80 ± 0 % % of [Privacy Breach Types] annually ... 0 % (x 0.0) 0 % (x 0.0) [Privacy Breach Types = Financ...] m 09 Jan 2020 ()-) Al Robustness 29.9 ± 0 10 ± 0 Δ: Status: 70 → Wish: 99.9 % #Intende∆%: $33 \pm 0\%$ $100 \pm 0\%$ % Intended Uptime for [Al Consumers] ... 0 % (x 0.0) 0% (x0.0) [Al Consumers = Transport,Auto...] 33% m 09 Jan 2025 ()-> Al Safety -90 ± 0 -42 ± 0 Δ: Status: 100 → Wish: 10 Number of ... 4%: 100 ± 0 % $47 \pm 0\%$ Number of annual Al Damage occurrenc.0 % (x 0.0) 0 % (x 0.0) [Al Accidents = Al Action Fail...] 47%

()-) Al Transparency Δ: Status: 10 > Wish: 99.998 % of select: % of selected [Al Cases] where we car ... (Al Cases - Decision for Machi 1

🛗 09 Jan 2022

88%

80%

My (TSG) Observations and Conclusions

- The most fundamental obstacle to AI Standards progress is WE NEED TO QUANTIFY AND STRUCTURE-**RICHLY** all **CRITICAL** STAKEHOLDER VALUES (ALL > 8)
- Next: we need much-deeper, morecomprehensive identification, and consensus, and detailed knowledge (their values and constraints) about AI/ AI STANDARDS **STAKEHOLDERS**





End Slide



brianmooredraws.com



SLIDES FOLDER

"In my day, kids didn't build massive, ransomwarespewing botnets. They got a paper route."

DOI to slides pdf: https://doi.org/10.13140/RG.2.2.10935.12960



