VERDI

# APPLYING THE HLEG'S ASSESSMENT LIST FOR TRUSTWORTHY AI (ALTAI) IN THE AUTOMOTIVE CONTEXT

David Andrews & Hristina Veljanova

virtual vehicle

Das Land Steiermark

Wissenschaft und Forschung, Gesundheit und Pflegemanagement

ZUKUNFTSFONDS STEIERMARK

UNI GRAZ

# Outlook

1. Assessment List for Trustworthy AI (ALTAI)

2. VERDI Criteria Catalogue for trustworthy highly-automated vehicles and driver assistance systems (SAE L3)

3. Cross-checking VERDI & ALTAI

4. Wrap up

# Imagine a world...



[Oppressive Silence 2016]

# ~~Imagine~~ Realise a world...



[Schwifty Memes 2019], see also Sorrel (2016)

# Assessment List for Trustworthy AI (ALTAI)

# AI HLEG & Trustworthy AI

**Principles**
- respect for human autonomy
- prevention of harm
- fairness
- explicability

**Key Requirements**
1. Human Agency and Oversight
2. Technical Robustness and Safety
3. Privacy and Data Governance
4. Transparency
5. Diversity, Non-discrimination and Fairness
6. Societal and Environmental Well-being
7. Accountability

**Assessment List**
- ALTAI

Trustworthy AI should be lawful, **ethical and robust**[1]

[1] AI HLEG (2019): "Ethics Guidelines for Trustworthy AI"

ALTAI Web Tool

# Screenshots ALTAI



[Own screenshot]

# Screenshots ALTAI



[Own screenshot]

# Screenshots ALTAI



[Own screenshot]

# ALTAI → VERDI

ALTAI: "is intended for **flexible use**: organisations can […] add elements to it as they see fit, taking into consideration the **sector** they operate in" [2]

ALTAI:
AI in general

→

VERDI:
Automotive context in specific

## Trustworthiness

[2] AI HLEG (2020): "The Assesment List For Trustworthy Artificial Intelligence (ALTAI)"

# VERDI Criteria Catalogue
## for trustworthy partially automated vehicles and driver assistance systems
## (SAE L3)

# Approach

**VERDI Criteria Catalogue**

| Beyond compliance | Values by design | Essential functions | Manufacturers as a target group |

# Methodology

TRUESSEC.eu

SCOTT
Secure Connected Trustable Things

Disciplinary Support Studies
(Ethics, Law, Sociology, Psychology)

AI HLEG's ALTAI

TRUESSEC.eu Criteria Catalogue for trustworthy ICT products and services

VERDI Criteria Catalogue

VERDI Core Areas → VERDI Criteria → VERDI Indicators

VERDI
Core Areas

Transparency

Privacy and good data governance

Fairness

Autonomy

Responsibility and accountability

Protection

| | |
|---|---|
| Transparency | The Core Area 'Transparency' encompasses provider's **information duties** towards the user regarding the **system's functionality and limitations** as well as the **data** that is processed by the system. Additionally, it also focuses on **information representation**. |
| Privacy and good data governance | This Core Area entails two aspects: (1) any **personal data** processed as part of the interaction with the system should be **protected**, and (2) the user should have the possibility to **contro**l that data. |
| Autonomy | Autonomy refers to the ADS providing the user with the possibility to **choose and make decisions** regarding the (non-)use of certain automation aspects and services as well as acknowledging other parties' **rights and freedoms**. |
| Fairness | Fairness stands for **preventing cases of discrimination** due to algorithmic biases and societal factors (e.g. socio-economic status) and considering effects and contributions towards **social in- and exclusion**. |
| Responsibility and accountability | **Respect** and **clear information** about the stipulation of **roles and liabilities**. It furthermore addresses the legitimate and reasonable **expectations** of the user and society in relation to the system's functionality and reliability. |
| Protection | This core area refers to the protection of users, other road users and the surrounding from any **harms and risks** that might be caused by the ADS, including physical harm (**safety**) and protection of software errors and data (**security**). |

# VERDI Criteria

| VERDI Core Areas | VERDI Criteria |
|---|---|
| 1. Transparency<br>2. Privacy and Good Data Governance<br>3. Autonomy<br>4. Fairness<br>5. Responsibility and Accountability<br>6. Protection<br><br>↓<br><br>**Trustworthiness** | **I.   Privacy**<br>  1. Minimised Collection, Processing and Use of Personal Data<br>  2. Transparent Processing of Personal Data<br>  3. Privacy Commitment<br><br>**II.  Communication**<br>  4. Information Representation<br>  5. Explainability<br>  6. Clear Stipulation of Roles and Duties<br><br>**III. Feedback Management and Dispute Resolution**<br>  7. Feedback and Complaint Management<br>  8. Ability to Redress<br>  9. Statement of Legal Compliance<br> 10. Appropriate Dispute Resolution<br><br>**IV. Protection**<br> 11. Established Oversight Mechanisms<br> 12. Secure Infrastructure<br> 13. Vehicle Safety<br><br>**V.  Fairness towards Society**<br> 14. Non-discrimination<br> 15. Avoiding Algorithmic Bias<br> 16. Social and Environmental Responsibility<br> 17. Open Data Approach |

VERDI

Criteria

## VERDI

Criteria & Indicators

## Information Representation

This criterion relates to how information is communicated to those interacting with the automated driving system directly or indirectly, which includes the driver and vehicle passengers as well as all other road users. It has the goal to ensure that the information is represented in a way that is user-friendly, relevant, easily accessible, visible, and free of charge.

### VERDI Indicators

1) Any information exchange or act of communication between the FRU/ driver and the ADS meets the following requirements. It is

   a) provided in a user-friendly manner, e.g.

      i) in a plain language (understandable to lay persons)

      ii) with the possibility to choose from several widely used languages

      iii) as long as necessary and as short as possible (depending on the situation and context)

   b) relevant to the context (no information overload)

   c) easily visible and accessible

2) ADS-relevant information is provided without extra costs.

3) Information about the currently operating level of automation is also given to other road users, while especially considering vulnerable road users, by using standardized ways of communication (e.g. audio signals or visible icons).

4) All kind of information is easily perceivable by elderly and persons with disabilities.

5) The ADS applies recent accessibility guidelines (e.g. from W3C in operation manuals, requirements related towards the vehicle users) to represent information.

# Cross-checking VERDI & HLEG's ALTAI

# Cross-checking VERDI & HLEG's ALTAI

- Structure
- Terminology
- Tool

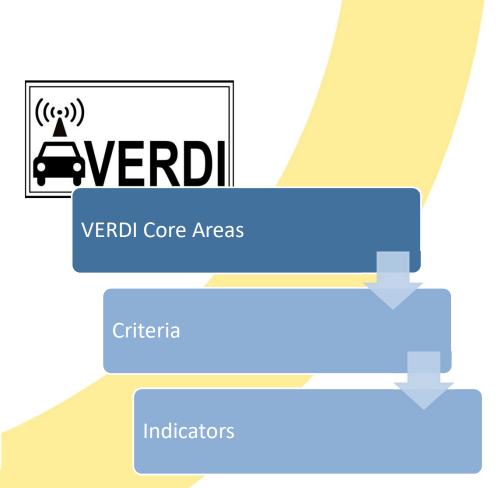INDEPENDENT
HIGH-LEVEL EXPERT GROUP ON
ARTIFICIAL INTELLIGENCE
SET UP BY THE EUROPEAN COMMISSION

AI

THE ASSESSMENT LIST FOR
TRUSTWORTHY ARTIFICIAL
INTELLIGENCE (ALTAI)
*for self assessment*

VERDI

[AI HLEG 2020]

# VERDI & HLEG's ALTAI: Structure

# VERDI Core Areas vs HLEG Ethical principles

| VERDI Core Areas | HLEG Ethical Principles | Cross-checking | |
|---|---|---|---|
| | | **VERDI** | **HLEG** |
| 1. Transparency | 1. Explicability | • Transparency | • Explicability |
| 2. Privacy and good data governance | 2. Respect for human autonomy | • Privacy and good data governance | • Respect for human autonomy<br>• Prevention of harm |
| 3. Autonomy | 3. Fairness | • Autonomy | • Respect for human autonomy |
| 4. Fairness | 4. Prevention of harm | • Fairness | • Fairness |
| 5. Responsibility and accountability | | • Responsibility and accountability | • Fairness |
| 6. Protection | | • Protection | • Prevention of harm |

# VERDI survey and next steps

## VERDI Criteria Assessment

You can find a short description of each VERDI core area by hovering over the respective core areas: Transparency, Privacy and good data governance, Autonomy, Fairness, Responsibility and accountability and Protection.

| 0% | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 10% | 20% | 30% | 40% | 50% | 60% | 70% | 80% | 90% | 100% |

### Information Representation

This criterion relates to how information is communicated to those interacting with the automated driving system directly or indirectly, which includes the driver and vehicle passengers as well as all other road users. It has the goal to ensure that the information is represented in a way that is user-friendly, relevant, easily accessible, visible, and free of charge.

| | 0 (None) | 1 (Low) | 2 (Medium) | 3 (High) | No answer |
|---|---|---|---|---|---|
| Transparency | ○ | ○ | ○ | ○ | ◉ |
| Privacy and good data governance | ○ | ○ | ○ | ○ | ◉ |
| Autonomy | ○ | ○ | ○ | ○ | ◉ |
| Fairness | ○ | ○ | ○ | ○ | ◉ |
| Responsibility and Accountability | ○ | ○ | ○ | ○ | ◉ |
| Protection | ○ | ○ | ○ | ○ | ◉ |

ⓘ Please rate to which extent the criterion addresses the corresponding core area from 0 (none) to 3 (high).

[Own screenshot]
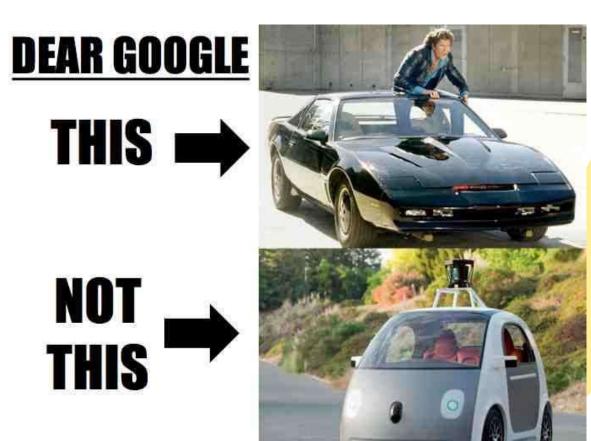
# Wrap up

https://verdi.uni-graz.at/en/

Final symposium in Spring 2021

Follow-up project on standardisation and certification

~~Imagine Realise~~ Shape a world...



[9GAG 2014]

# References

AI HLEG (2019): Ethics Guidelines for Trustworthy AI. [online] https://ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai [24.02.2021]

AI HLEG (2020): The Assesment List For Trustworthy Artificial Intelligence (ALTAI). [online] https://ec.europa.eu/digital-single-market/en/news/assessment-list-trustworthy-artificial-intelligence-altai-self-assessment [24.02.2021].

Sorrel, Charlie (2016): Self-Driving Mercedes Will Be Programmed To Sacrifice Pedestrians To Save The Driver. In: Fastcompany.com [online] https://www.fastcompany.com/3064539/self-driving-mercedes-will-be-programmed-to-sacrifice-pedestrians-to-save-the-driver [24.02.2021]

**Images**

9GAG (2014): Google's new self driving car is unacceptable…[online] https://9gag.com/gag/aYb5ZeN [24.02.2021]

Oppressive Silence (2016): self driving car. [online] https://imgur.com/gallery/I0ivc64 [24.02.2021]

Schwifty Memes (2019): [online] https://schwifty-memes.tumblr.com/post/189712495869 [24.02.2021]

# Thank you for your attention!

david.andrews@uni-graz.at
hristina.veljanova@uni-graz.at

*We work for*
**tomorrow**

UNI
GRAZ