



SIMO KYLLÖNEN

---

# AVOIN TIEDE JA TEKÖÄLY

EETTISET HAASTEET JA HYVÄT KÄYTÄNNÖT

The open science approach, when combined with AI technologies, can significantly contribute to addressing the current pressing global threats.

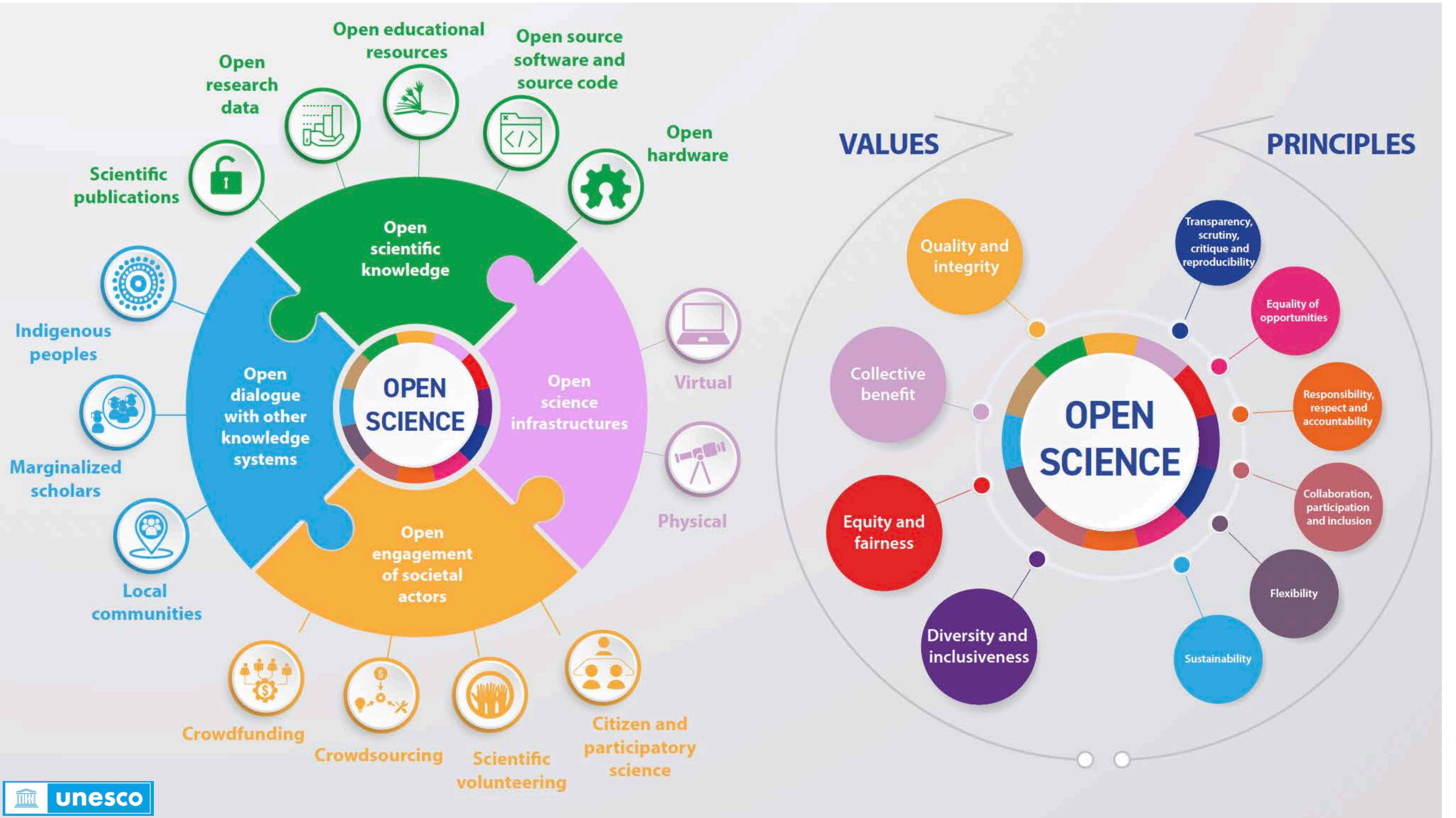
## **NASA Fights *Climate Change* with Open Science and AI**

### **Building a Powerful AI Tool for Earth Science**

In 2023, NASA teamed up with IBM to create a powerful AI tool for studying Earth. This tool is like a basic building block, trained on a massive dataset of NASA's satellite images. Anyone can use this tool for free to build their own studies on climate challenges.

Think of it like a Swiss Army Knife for science. It can't do everything by itself, but it can be a starting point for lots of different projects. Scientists can take this basic tool and add specific information to make it even better at a particular task.





# Open Science at the Generative AI Turn: An Exploratory Analysis of Challenges and Opportunities

Mohammad Hosseini\*<sup>1</sup>, Serge P.J.M. Horbach<sup>2</sup>, Kristi Holmes<sup>1,3</sup>, Tony Ross-Hellauer<sup>4</sup>

<sup>1</sup> Department of Preventive Medicine, Northwestern University Feinberg School of Medicine, Chicago, Illinois, United States of America

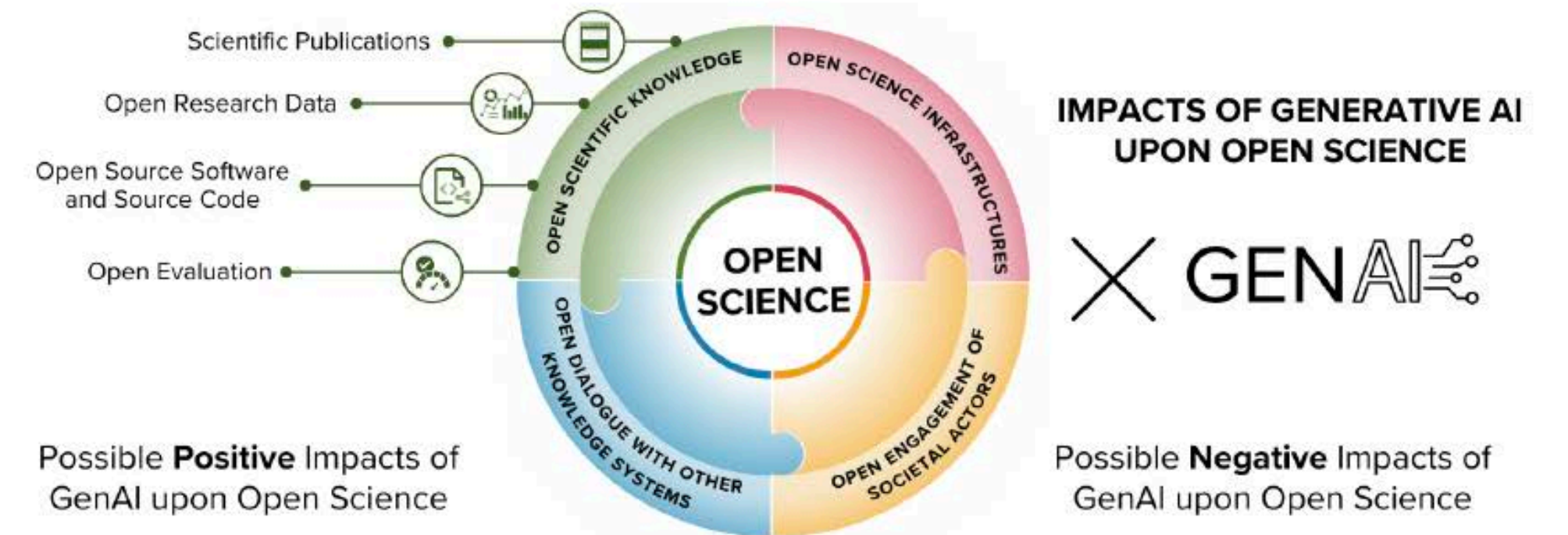
<sup>2</sup> Institute for Science in Society, Radboud University, Nijmegen, The Netherlands

<sup>3</sup> Galter Health Sciences Library and Learning Center, Northwestern University Feinberg School of Medicine, Chicago, Illinois, United States of America

<sup>4</sup> Open and Reproducible Research Group, Know-Center GmbH and Institute for Interactive Systems and Data Science, Graz University of Technology, Graz, Austria

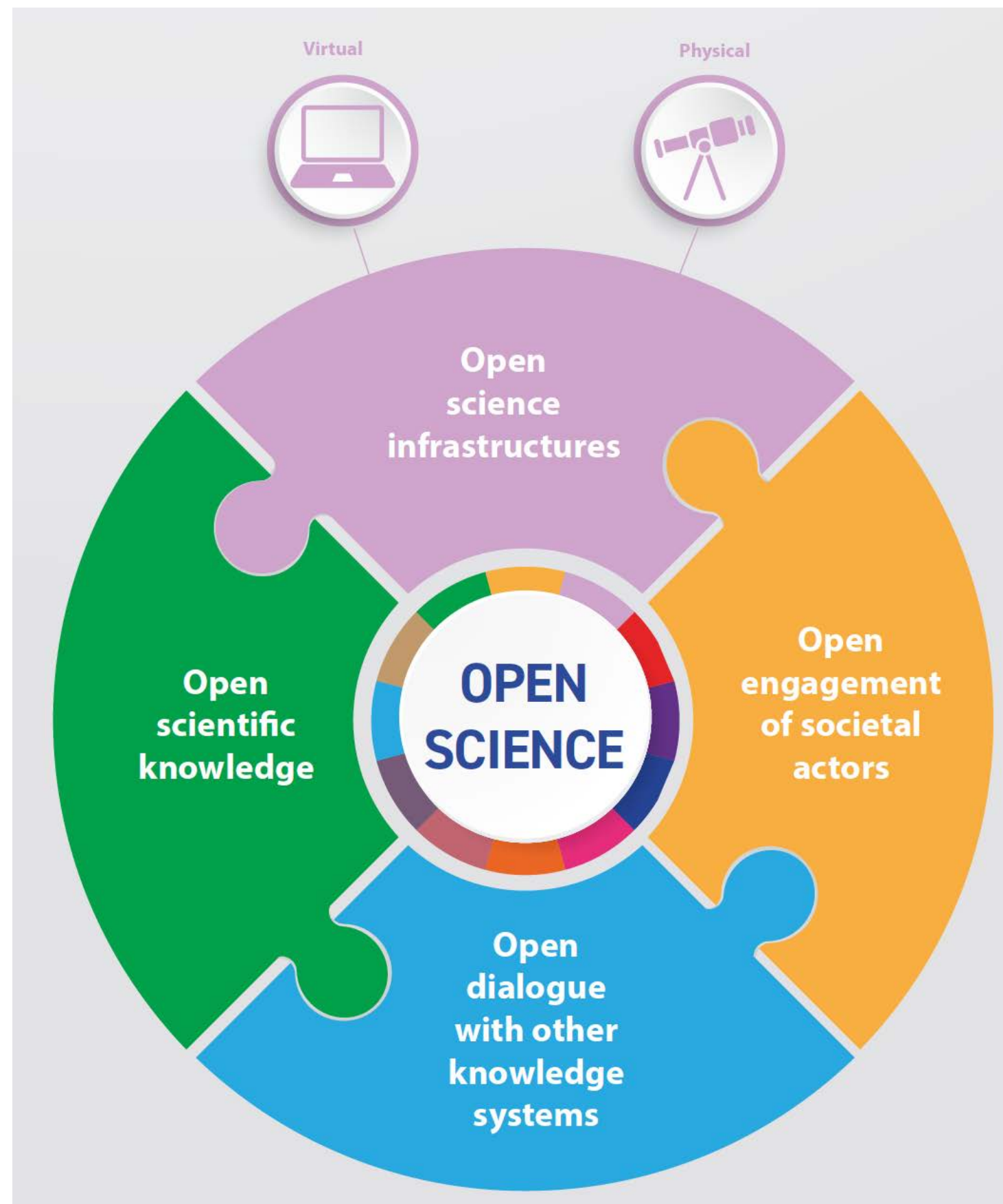
\* Corresponding author

Email: mohammad.hosseini@northwestern.edu



Open Scientific Knowledge	
<ul style="list-style-type: none"> <li>Transitioning from open to meaningful and equitable access for researchers, policymakers and laypersons.</li> <li>Helping researchers cope with information overload.</li> </ul>	<ul style="list-style-type: none"> <li>Enabling paper mills and fake publications, and increased biases and misinformation about science.</li> <li>Deterring researchers from open sharing of research out of a fear of having content harvested without attribution.</li> </ul>
<ul style="list-style-type: none"> <li>Supporting Research Data Management and data sharing documentation, data analysis and curation.</li> <li>Detecting errors or inconsistencies in datasets.</li> <li>Creating synthetic datasets.</li> </ul>	<ul style="list-style-type: none"> <li>Generating fabricated datasets that support specific hypotheses and falsely validate fake papers.</li> <li>Worsening reproducibility and integrity by enabling the use of unvalidated mathematical and statistical methods.</li> </ul>
<ul style="list-style-type: none"> <li>Increasing productivity and efficiency in coding.</li> <li>Reviewing code in real-time, thereby increasing quality of publicly available code.</li> <li>Improving code documentation and reusability.</li> </ul>	<ul style="list-style-type: none"> <li>Increasing code inaccuracies (due to indeterminism and variability of outputs) and biases.</li> </ul>
<ul style="list-style-type: none"> <li>Enhancing evaluation of scholarly contributions.</li> <li>Increased reviewer pool and lower burden on reviewers.</li> </ul>	<ul style="list-style-type: none"> <li>Increasing the risks of fake review reports, leading to reduced value of open reports.</li> <li>Compromising confidentiality.</li> </ul>
Open Science Infrastructures	
<ul style="list-style-type: none"> <li>Enhancing efficient use of open repositories and scholarly indices.</li> </ul>	<ul style="list-style-type: none"> <li>Compromising reproducibility since most major models are not themselves open.</li> <li>Lacking proper attribution and risk of infringing copyrights or committing plagiarism.</li> <li>Risking monetization of research information and leading to further restrictions.</li> </ul>
Open Engagement of Societal Actors	
<ul style="list-style-type: none"> <li>Increasing opportunities for effective engagement of more parties.</li> <li>Leveling the playing field for citizen scientists.</li> </ul>	<ul style="list-style-type: none"> <li>Increasing the likelihood of technological dependence and inequity.</li> <li>Spreading misinformation and increasing the risks of shaping public opinion based on false information about research.</li> </ul>
Open Dialogue with Other Knowledge Systems	
<ul style="list-style-type: none"> <li>Facilitating knowledge exchange across different epistemic communities.</li> <li>Enhancing translation of open outputs and enabling dialogue.</li> </ul>	<ul style="list-style-type: none"> <li>Increasing biases that reflect global hegemonies and serve the interests of specific groups.</li> <li>Increasing the likelihood of mistranslation and misrepresentations, especially of those languages less reflected in training data.</li> </ul>

## AVOIMEN TIETEEN INFRA



- ▶ **AI** voi tehostaa avointen aineistojen ja arkistojen käyttöä
- ▶ **Mutta:**
  - ▶ Käyttöä rajoittaa AI:n läpinäkyvyyteen liittyvät riskit
    - ▶ LLMt ym eivät itse ole aidosti avoimia
    - ▶ plagioinnin ja tekijänoikeuksien rikkominen
    - ▶ aineistojen taloudellinen hyödyntäminen

# Opening up ChatGPT: tracking openness of instruction-tuned LLMs

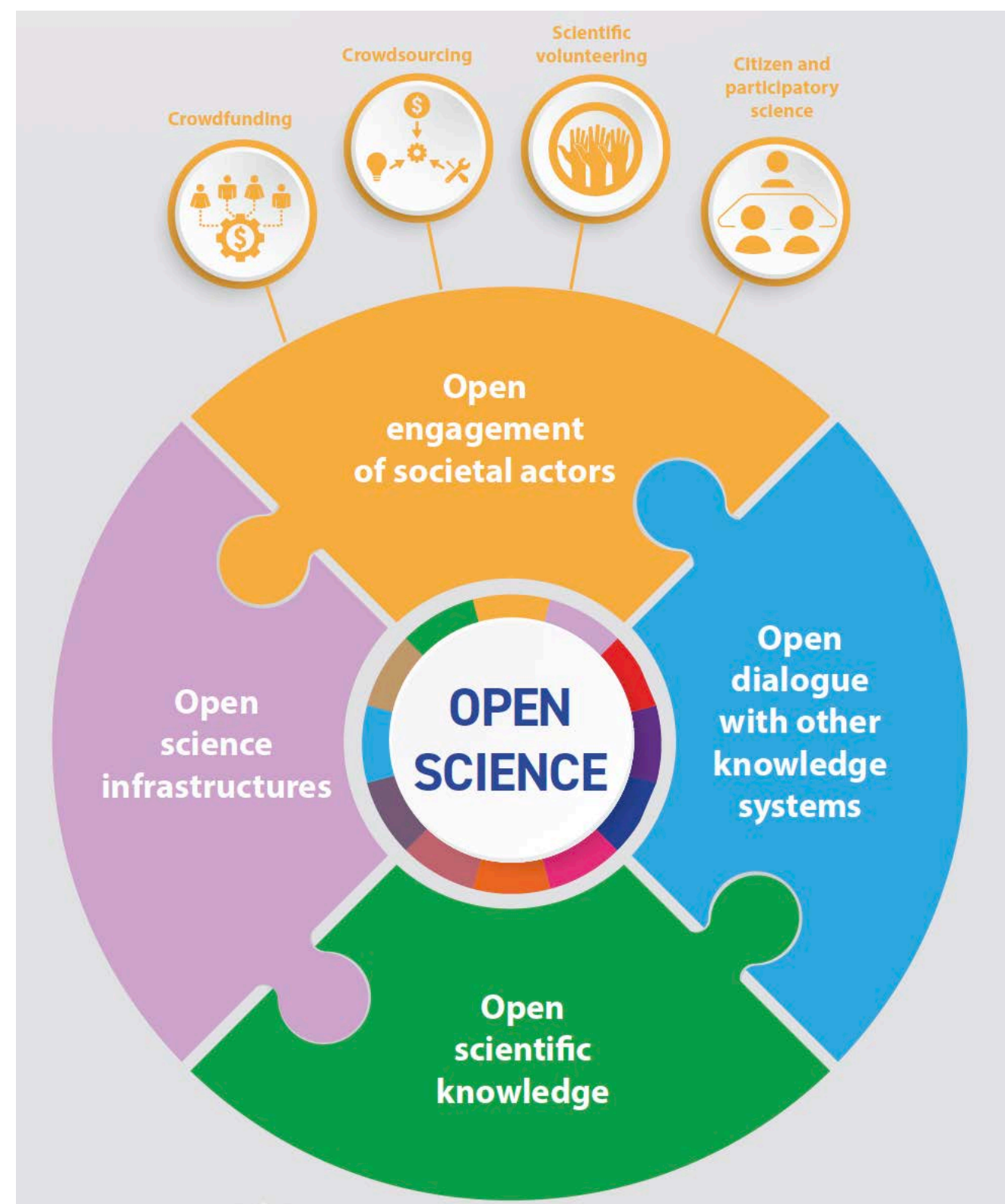
⚡ **FAccT'24 paper** ⚡ Liesenfeld, Andreas, and Mark Dingemans. 2024. 'Rethinking Open Source Generative AI: Open-Washing and the EU AI Act'. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency (FAccT '24)*. Rio de Janeiro, Brazil: ACM. ([PDF](#)).

There is a growing amount of instruction-tuned text generators billing themselves as 'open source'. How open are they really? [FAccT'24](#) [CUI'23](#) [repo](#)

Project (maker, bases, URL)	Availability					Documentation						Access		
	Open code	LLM data	LLM weights	RL data	RL weights	License	Code	Architecture	Preprint	Paper	Modelcard	Datasheet	Package	API
OLMo 7B Instruct Ai2	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✓	✓	✓	~
BLOOMZ bigscience-workshop	✓	✓	✓	✓	~	~	✓	✓	✓	✓	✓	✓	✗	✓
AmberChat LLM360	✓	✓	✓	✓	✓	✓	~	~	✓	✗	~	~	✗	✓
Open Assistant LAION-AI	✓	✓	✓	✓	✗	✓	✓	✓	~	✗	✗	✗	✓	✓
OpenChat 3.5 7B Tshinghua University	✓	✗	✓	✗	~	~	~	~	~	✗	~	~	✗	✗
Pythia-Chat-Base-7... togethercomputer	✓	✓	✓	✓	~	~	~	~	~	✗	~	~	✗	✗
Cerebras GPT 111... Cerebras + Schramm	~	✓	✓	✓	~	~	~	~	~	✗	~	~	✗	✗
RedPajama-INCITE... TogetherComputer	~	✓	✓	✓	~	~	~	~	~	✗	~	~	✗	✗

Falcon-180B-chat Technology Innovation In...	✗	~	~	~	~	✗	✗	~	~	✗	~	~	✗	✗	✗
Gemma 7B Instruct Google DeepMind	~	✗	~	✗	~	✗	✗	~	~	✗	~	~	✗	✗	✗
Orca 2 Microsoft Research	✗	✗	~	✗	✓	✗	✗	~	~	✗	~	~	✗	✗	~
Command R+ Cohere AI	✗	✗	✗	✓	✓	~	✗	✗	✗	✗	✗	✗	~	✗	✗
LLaMA2 Chat Facebook Research	✗	✗	~	✗	~	✗	✗	~	~	✗	~	~	✗	✗	~
Nanbeige2-Chat Nanbeige LLM lab	✓	✗	✗	✗	✓	~	✗	✗	✗	✗	✗	✗	✗	✗	~
Llama 3 Instruct Facebook Research	✗	✗	~	✗	~	✗	✗	~	~	✗	✗	✗	~	✗	~
Solar 70B Upstage AI	✗	✗	~	✗	~	✗	✗	~	~	✗	✗	✗	~	✗	~
Xwin-LM Xwin-LM	✗	✗	~	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	~
ChatGPT OpenAI	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	~	✗	✗	✗	✗

## AVOIN TUTKIMUSKULTTUURI JA OSALLISTUMINEN



- ▶ AI voi madaltaa kynnystä osallistua ja auttaa tavoittamaan laajemman yleisön

- ▶ **Mutta:**

- ▶ Teknologiaan liittyvien epäoikeudenmukaisuuksien lisääntyminen
- ▶ Virheellisen informaation levittämiseen ja manipulatiiviseen käyttöön liittyvät riskit

- Discover
- Submit
- Welcome to the AIID
- Discover Incidents
- Spatial View
- Table View
- List view
- Entities
- Taxonomies
- Submit Incident Reports
- Submission Leaderboard
- Blog
- AI News Digest
- Risk Checklists
- Random Incident
- Sign Up

Incidents
Issue Reports
Reports

Show Live data
 Reset filters

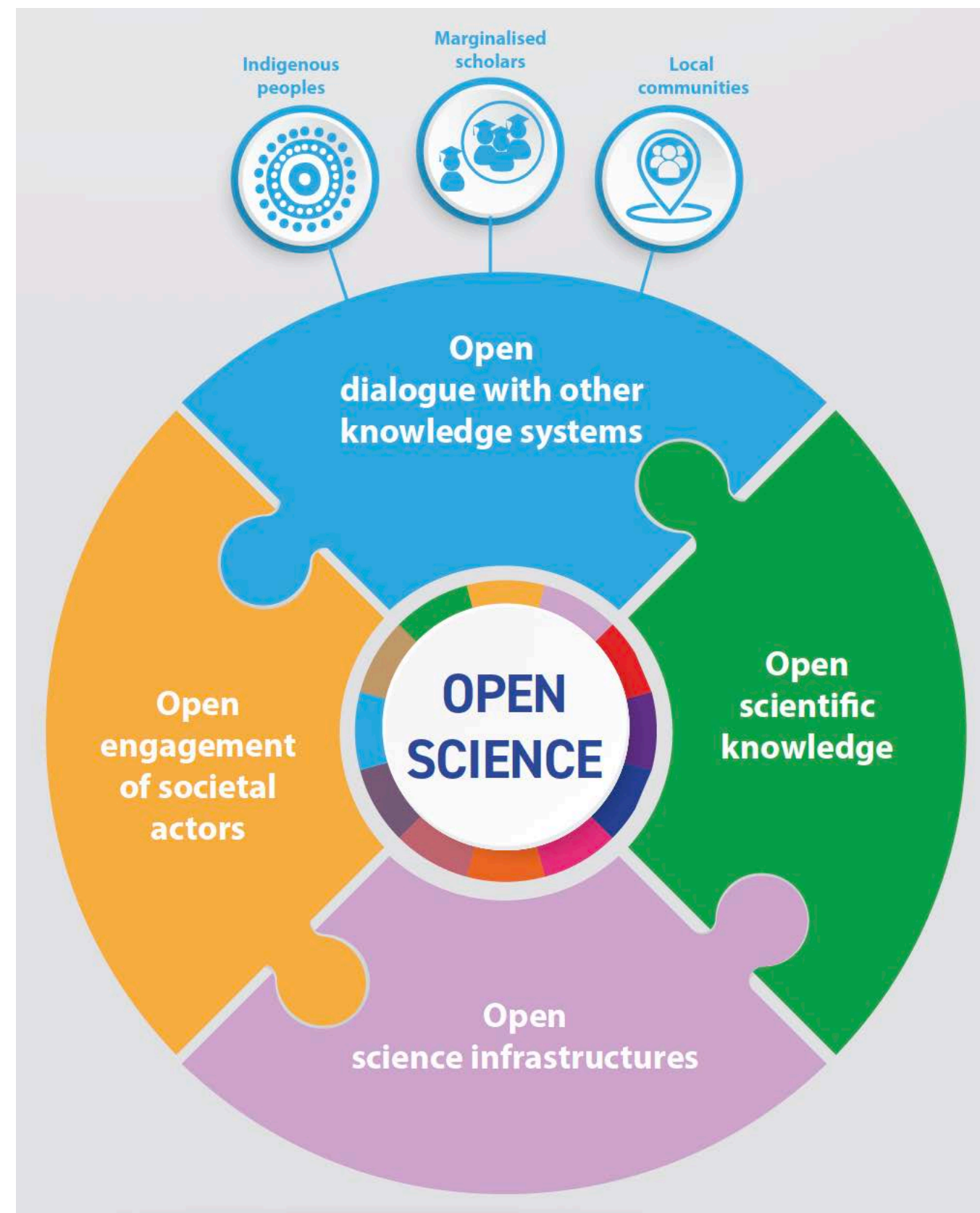
Displaying 10 of 836 incidents

INCIDENT ID	TITLE	DESCRIPTION	DATE	ALLEGED DEPLOYER OF AI SYSTEM	ALLEGED DEVELOPER OF AI SYSTEM	ALLEGED HARMED OR NEARLY HARMED PARTIES
<a href="#">Incident 1</a>	Google's YouTube Kids App Presents Inappropriate Content	YouTube's content filtering and recommendation algorithms exposed children to disturbing and inappropriate videos.	2015-05-19	<a href="#">YouTube</a>	<a href="#">YouTube</a>	<a href="#">Children</a>
<a href="#">Incident 23</a>	Las Vegas Self-Driving Bus Involved in Accident	A self-driving public shuttle by Keolis North America and Navya was involved in a collision with a human-driven delivery truck in Las Vegas, Nevada on its first day of service.	2017-11-08	<a href="#">Navya, Keolis North America</a>	<a href="#">Navya, Keolis North America</a>	<a href="#">Navya, Keolis North America, bus passengers</a>
<a href="#">Incident 4</a>	Uber AV Killed Pedestrian in Arizona	An Uber autonomous vehicle (AV) in autonomous mode struck and killed a pedestrian in Tempe, Arizona.	2018-03-18	<a href="#">Uber</a>	<a href="#">Uber</a>	<a href="#">Elaine Herzberg, pedestrians</a>
<a href="#">Incident 12</a>	Common Biases of Vector Embeddings	Researchers from Boston University and Microsoft Research, New England demonstrated gender bias in the most common techniques used to embed words for natural language processing (NLP).	2016-07-21	<a href="#">Microsoft Research, Boston University</a>	<a href="#">Microsoft Research, Google, Boston University</a>	<a href="#">Women, Minority Groups</a>
<a href="#">Incident 5</a>	Collection of Robotic Surgery Malfunctions	Study on database reports of robotic surgery malfunctions (8,061), including those ending in injury (1,391) and death (144), between 2000 and 2013.	2015-07-13	<a href="#">Hospitals, Doctors</a>	<a href="#">Intuitive Surgical</a>	<a href="#">patients</a>
<a href="#">Incident 6</a>	TayBot	Microsoft's Tay, an artificially intelligent chatbot, was released on March 23, 2016 and removed within 24 hours due to multiple racist, sexist, and anti-semitic tweets generated by the bot.	2016-03-24	<a href="#">Microsoft</a>	<a href="#">Microsoft</a>	<a href="#">Twitter Users</a>
<a href="#">Incident 10</a>	Kronos Scheduling Algorithm Allegedly Caused Financial Issues for Starbucks Employees	Kronos's scheduling algorithm and its use by Starbucks managers allegedly negatively impacted financial and scheduling stability for Starbucks employees, which disadvantaged wage workers.	2014-08-14	<a href="#">Starbucks</a>	<a href="#">Kronos</a>	<a href="#">Starbucks employees</a>
<a href="#">Incident 11</a>	Northpointe Risk Models	An algorithm developed by Northpointe and used in the penal system is two times more likely to incorrectly label a black person as a high-risk re-offender and is two times more likely to incorrectly label a white person as low-risk for reoffense according to a ProPublica	2016-05-23	<a href="#">Northpointe</a>	<a href="#">Northpointe</a>	<a href="#">Accused People</a>

<https://incidentdatabase.ai/apps/incidents/?view=incidents>



## AVOIN DIALOGI



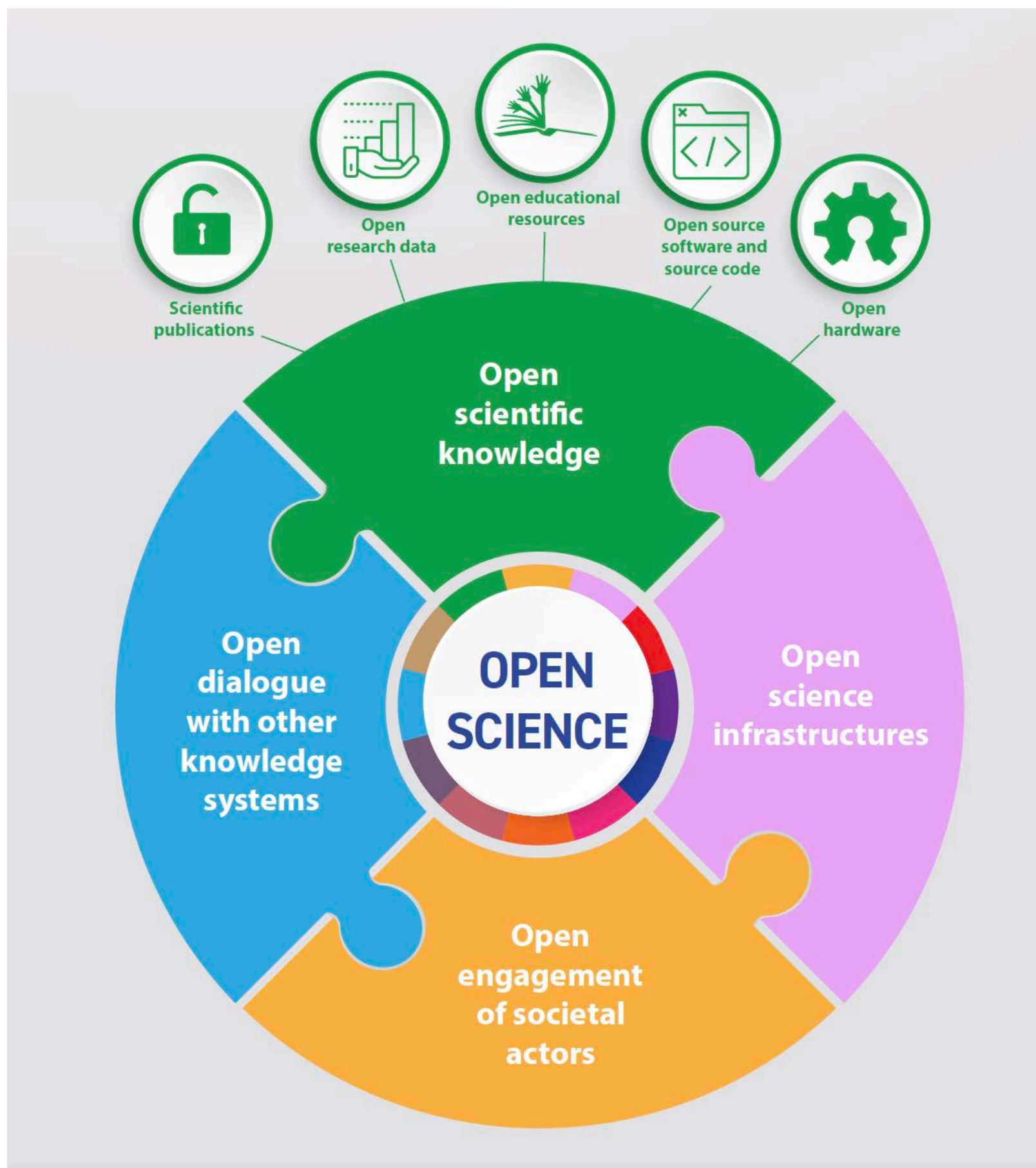
- ▶ AI voi helpottaa dialogia muiden episteemisten yhteisöjen kanssa:

- ▶ Avointen tutkimustuotosten kääntäminen

- ▶ **Mutta:**

- ▶ Lisääntynyt painotus valtakulttuurin avoimena saataviin näkemyksiin, oletuksiin
  - ▶ Käännös- ja tulkintavirheiden lisääntyminen

# TUTKIMUSAINEISTOJEN, -MENETELMIEN JA -JULKAISUJEN AVOIN SAATAVUUS



- ▶ Aito, **mielekäs** ja **tasapuolinen** saatavuus
    - ▶ tieteellisen tiedon **lukutaito**
    - ▶ **informaatioähkyn** hallinta
  - ▶ **AI** voi edistää näitä tarjoamalla tietoa, tuloksia ymmärrettävällä, kootulla ja kuratoidulla tavalla
- ▶ **Mutta:**
- ▶ **AI** helpottaa artikkelitehtailua, vale/vinoutuneen informaation leviämistä
  - ▶ **AI** voi rajoittaa tutkijoita avaamasta tutkimustaan ja julkaisujaan

“Ensimmäinen kokemukseni oli viime kesäkuussa. Olin arvioimassa konferenssiin submittoitua abstraktia, joka oli melko helppo tunnistaa tekoälyn tuottamaksi. Siinä ei periaatteessa ollut mitään vikaa, mutta jos on vähänkin sisällä keskustelussa, huomaa, että kyseessä on tilkkutäkki eri tasoista keskusteluista: johdantokirjan kaltaista tekstiä yhdistetään huippuspesialisoituneeseen keskusteluun ilman, että kokonaisuuden taustalta kuuluu selkeää ääni ja suunta.

Toinen kokemus oli rajumpi, ja uskallan väittää, että oli kyse onnekkaasta sattumasta - siitä, että juuri minä satuin olemaan ko. käsikirjoituksen vertaisarvioija - että käsikirjoitus hylättiin. Kyse oli tutkimaani aihetta käsittelevästä artikkelikäsikirjoituksesta, joka oli submittoitu alani JUF0 1-tasoiseen kv-lehteen. Pintapuolisesti luettuna jutussa ei ollut mitään erityistä vikaa, olkoonkin, että argumentti ei ollut erityisen selkeästi jäsentynyt ja yhdisteli teemoja, jotka pulpahtavat esiin googlaamalla ko. aiheesta. Mutta tällaisia kässäreitä nyt on liikkeellä paljonkin, ja usein (etenkin kokemattomammilla) refereetit antavat niille mahdollisuuden eli R&R.

Olen siitä 98 % varma, että kyse oli AI:n tuottamasta tekstistä. Näin siksi, että käsikirjoituksessa esitettiin yksi minun julkaisemani väite, josta tiedän varmuudella, että kukaan muu ei ole sitä esittänyt. Väitettä ei myöskään avattu, vaan se esitettiin triviaalina seikkana (mitä se ei totisesti ole). Samoin triviaaleina esitettiin kahden muun tutkijan julkaisemat spesifit väitteet, joista toisen premissit ovat suorassa ristiriidassa mainitun oman väitteeni kanssa, mutta tavalla, jota kone ei (vielä) tunnista.

Nämä aihetodisteet riittivät minulle perusteeksi esittää kässärin hylkäämistä ja kerroin päätoimittajalle luottamuksellisesti epäileväni AI:ta.”

“Oma analyysini on, että nykyinen AI pystyy tuottamaan uskottavan tuntuista tekstiä, joka hyvinkin saattaisi mennä läpi ainakin heikommassa lehdissä. Sen heikkous on toistaiseksi, että aidon argumentaation ja pohdinnan sijaan tyyli on lähempänä oppikirjamaista tekstiä: väitteet, joihin alansa tuntevat tutkijat suhtautuisivat kriittisesti ja kyseenalaistaen, esitetään ikään kuin faktoina.

**Kauhuskenaarioni** toinen vaihe on mielestäni se, että AI:lle esitetään pyyntö muokata tyyliä enemmän “alan tyylin mukaiseksi”. TAI sille syötetään esim. referee-lausunto, jossa yksityiskohtaisesti selostetaan, mikä argumentaatiossa on pielessä. Kun AI sitten toteuttaa nämä toiveet, on AI-lähtöisyyden tunnistamisen entistä vaikeampaa. Pelkään, että AI kehittyy tässä suhteessa kovaa vauhtia (ilmainen versio on selvästi huonompi kuin maksullinen, nykyinen ilmainen versio on todella paljon parempi kuin vuoden takainen, jne.). Tästä syystä onkin huolestuttavaa, että esim. Taylor & Francis myy korpuksensa miljoonadiilillä Microsoftin AI:n kehittämistarkoitukseen (kertomatta asiasta kirjoittajille). Koska kaikki OA-materiaali on materiaalia AI:lle, OA:n edistäminen syöttää suoraan lapaan.

Kauhuskenaarion kolmas vaihe (jossa uskoakseni saatamme jo olla) on se, että materiaali, johon AI perustuu, on osittain AI:n itsensä tuottamaa. Tämä toimijuutta - intentiota ja vastuuta - vailla oleva teksti sitten leviää hallitsemattomasti, viime kädessä myös ihmiskirjoittajiin.”



Tutkijoiden halu julkaista avoimesti vähenee

Avoimen arvioinnin edellytykset heikkenevät

Riski vinoutuneen/valetiedon leviämiseen kasvaa

▸ **vastoin avoimen tieteen tavoitteita:**

▸ **tasapuolisuus, kollektiivinen hyöty, diversiteetti, läpinäkyvyys**

# TEKOÄLY HAASTAVA TOIMINTAYMPÄRISTÖ AVOIMELLE TIETEELLE

TEKOÄLYAVUSTEISET TEKNIIKAT



AVOIN TIEDE



JA BISNESLOGIIKAT,  
TOIMINTAKULTTUURIT ...

“Because generative AI has been explosive and organizations often need help to establish basic OS policies, these concerns have not yet reached the most visible parts of the OS ecosystem.

... We are concerned that if OS policymakers continue to not explicitly consider AI in their debates, AI harm will continue to propagate and policy creation to prevent it will become increasingly difficult.”

Acion, Laura et al. (2023) *Nature*

# ONKO AVOIN TIEDE VAIVANSA VÄÄRTI?

## HYVÄT KÄYTÄNNÖT?

Ohjeet ja  
suositukset

Avoimen tieteen, tekoälyn ja  
tutkimusetiikan integrointi ohjeissa

Vastuusta välittävä yhteisö

Yhteisön tuki tekoälyn vastuulliselle  
käytölle tutkimuksen eri  
vaiheissa ja julkaisuissa:  
koulutus, palvelut,  
palkitseminen

Vastuullisuutta helpottavat välineet

Laadukkaat avoimet julkaisut ja tietoturvalliset  
tekoälyvälineet tutkijoille (esim. Copilot, CurreChat)

Vastuullisuuden mahdollistavat rakenteet

Tasapuolisuutta, diversiteettiä ja tieteen  
integriteettiä ja tutkijoiden oikeuksia turvaavat  
avoimen tieteen ja tekoälyn rakenteet

## GENERATIIVISEN TEKOÄLYN KÄYTTÖ TUTKIMUKSESSA

Helsingin yliopisto tukee generatiivisen tekoälyn vastuullista ja harkittua hyödyntämistä tutkimuksessa. Tämä sivu sisältää tutkijalle tarkoitettuja ohjeita generatiiviseen tekoälyyn perustuvien työkalujen käyttöön.

Generatiivisen tekoälyn käyttöä tutkimuksessa koskeva yliopiston linjaus perustuu Euroopan tutkimusalueen julkaisemaan [ohjeistukseen generatiivisen tekoälyn vastuullisesta käytöstä tutkimuksessa](#) (sivustolta suomeksi saatavilla konekäännös).

Helsingin yliopiston tutkijoiden on noudatettava kaikessa toiminnassaan yliopiston julkaisemia [tekoälyn käytön yleisiä periaatteita](#) (Linkki vie Flammaan, joka vaatii kirjautumisen). Pehdy yleisiin periaatteisiin ensin.

Tutkimusta koskeva tekoälyohjeistus täydentää [Tutkimuseettisen neuvottelukunnan \(TENK\) ohjeita hyvästä tieteellisestä käytännöstä](#). Generatiivinen tekoäly on vain työkalu – tutkimuksen vakiintuneet eettiset periaatteet koskevat myös tekoälyn tuella tehtävää tutkimusta, ja tekoälytyökalujen vastuullisessa ja eettisessä käytössä pääsee pitkälle terveellä järjellä.

Helsingin yliopiston tietotekniikkakeskus on koonnut käytännön ohjeita [generatiivisen tekoälyn käytöstä yliopistolla](#). Ohjeet sisältävät vinkkejä erilaisten sovellusten, kuten [Copilot](#) ja [CurreChat](#) (ChatGPT:n pohjalta yliopistolaisille räätälöity työkalu), sekä [Kontra](#) (konekäännin, joka kääntää englannin, ruotsin ja suomen välillä) käyttöön.

Tietotekniikkakeskuksen käytännön ohjeet antavat myös suuntaviivoja ulkoisten tekoälysovellusten, kuten ChatGPT:n ja kuvagenerointipalvelu Midjourneyn käyttöön. Ajan myötä yliopistolaisten käytössä olevien uusien vastaavien työkalujen määrä todennäköisesti kasvaa.

Tämä ohjeistus koskee generatiivisen tekoälyn **käyttöä** tutkimuksessa, ja se kattaa kaikki generatiivisen tekoälyn tyypit. Ohjetta päivitetään tarvittaessa. Tulevaisuudessa yliopisto harkitsee ohjeistuksen täydentämistä generatiivisten tekoälysovellusten kehittämistä koskevalla ohjeella.

Sisällysluettelo

Avaa

# Kiitos!

Simo Kyllönen

Yliopistonlehtori

Helsingin yliopisto

<https://researchportal.helsinki.fi/en/persons/simo-kyllonen>

[simo.kyllonen@helsinki.fi](mailto:simo.kyllonen@helsinki.fi)

© 2023 by Simo Kyllönen is licensed under CC BY 4.0