

# Ориентиране в етика на ИИ

---

Образователен форум „Образованието през  
парадигмата на хуманизма и изкуствения интелект“,  
Варна

---

ДИМИТЪР НИКОЛОВ, Д-Р  
d\_nikolov@tu-sofia.bg

ТЕХНИЧЕСКИ УНИВЕРСИТЕТ СОФИЯ  
ФЕТТ

АДРЕС  
бул. Климент Охридски 8, гр. София

# За мен



## доц. д-р Димитър Николов

- ФЕТТ при ТУ-София
  - Микро Нано Лаборатория, СТП АД - ръководител
  - Клъстер Микроелектроника и индустриални електронни системи - председател на УС
  - AI клъстер България
  - 50+ публикации
  - 15 Европейски проекти
- [d\\_nikolov@tu-sofia.bg](mailto:d_nikolov@tu-sofia.bg)

# Съдържание

03

Въведение

04

Propublica <->  
COMPASS

05

Етични норми и  
тяхното влияние 1

06

Етични норми и  
тяхното влияние 2

07

Законодателство  
2016 - 2023  
географски

08

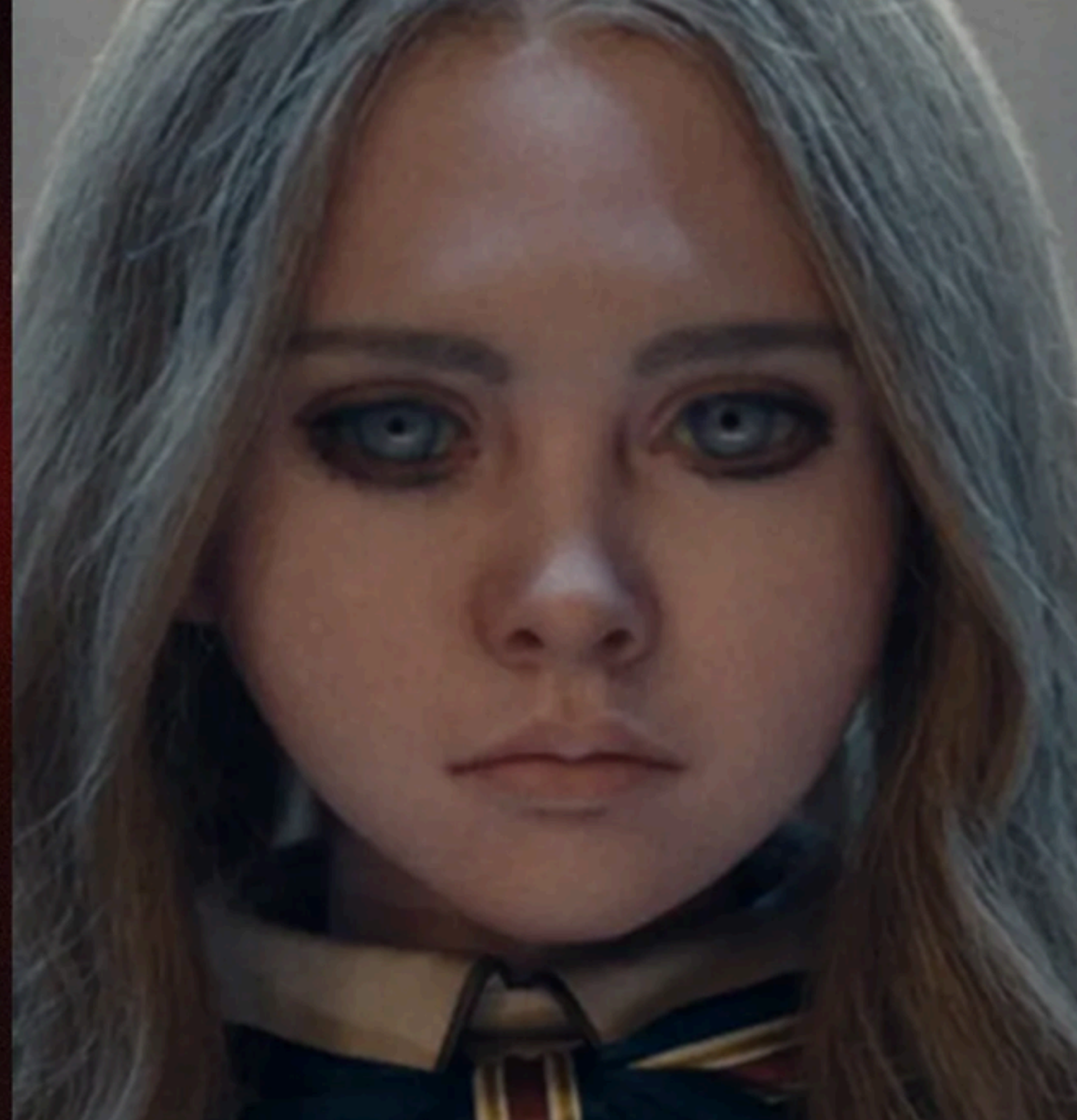
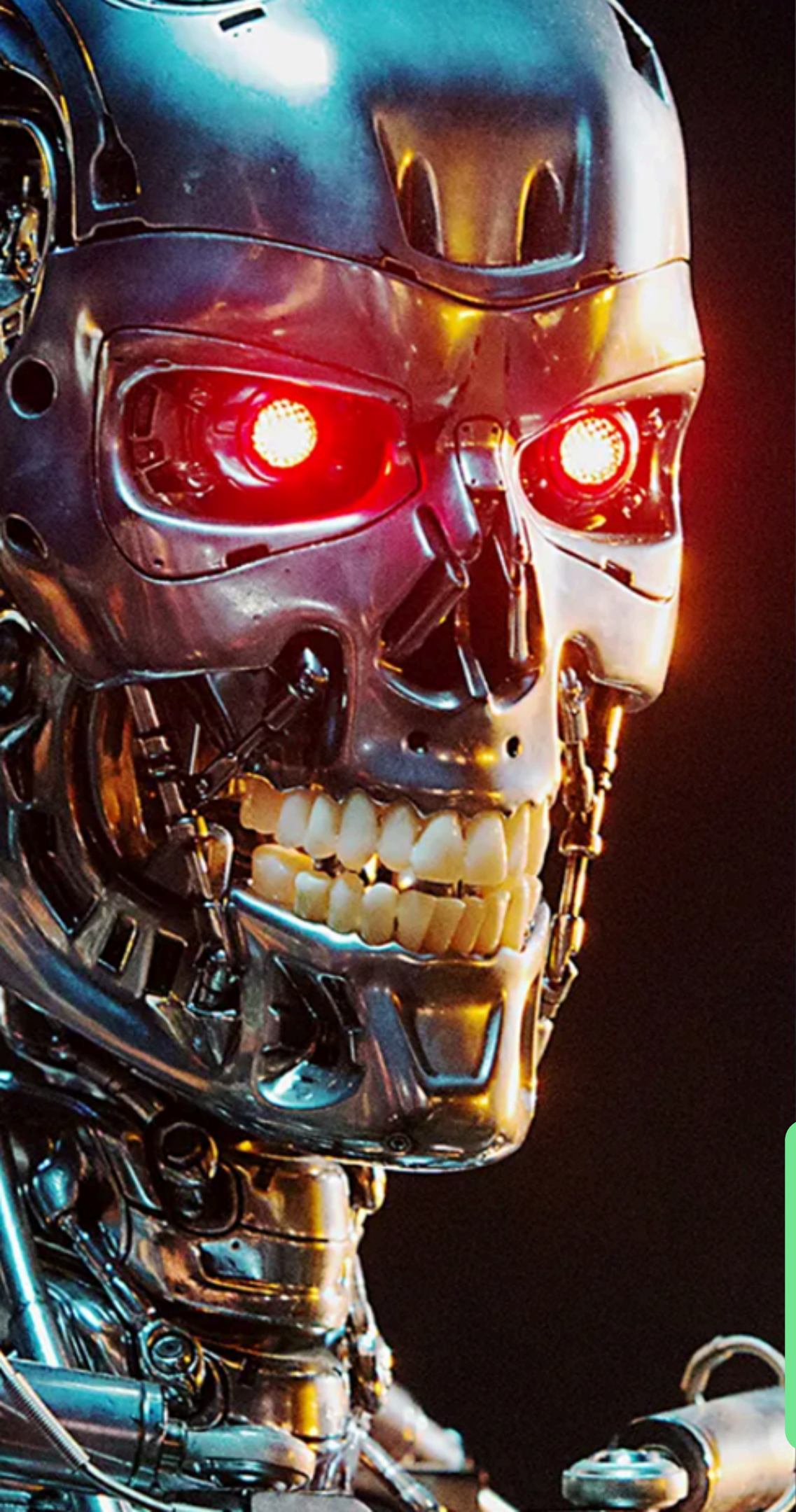
Законодателство  
2016 - 2023 по  
области

09

EU AI Act

10

Z-inspection  
process



**"Невъзможно е да се състави набор от правила, претендиращи да опишат какво човек трябва да прави във всякакви възможни обстоятелства."**

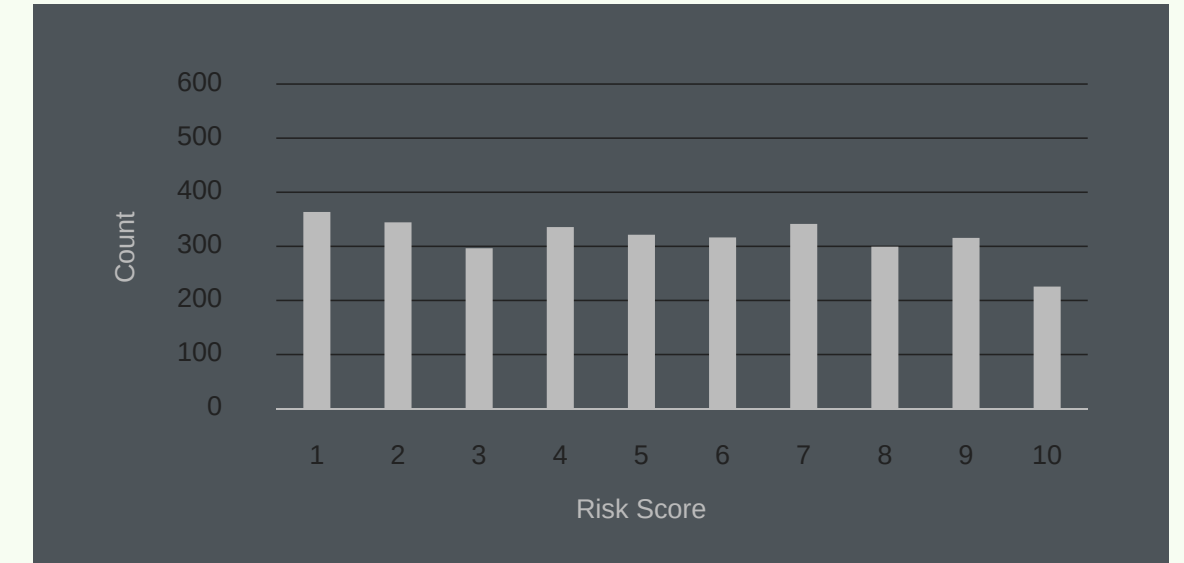
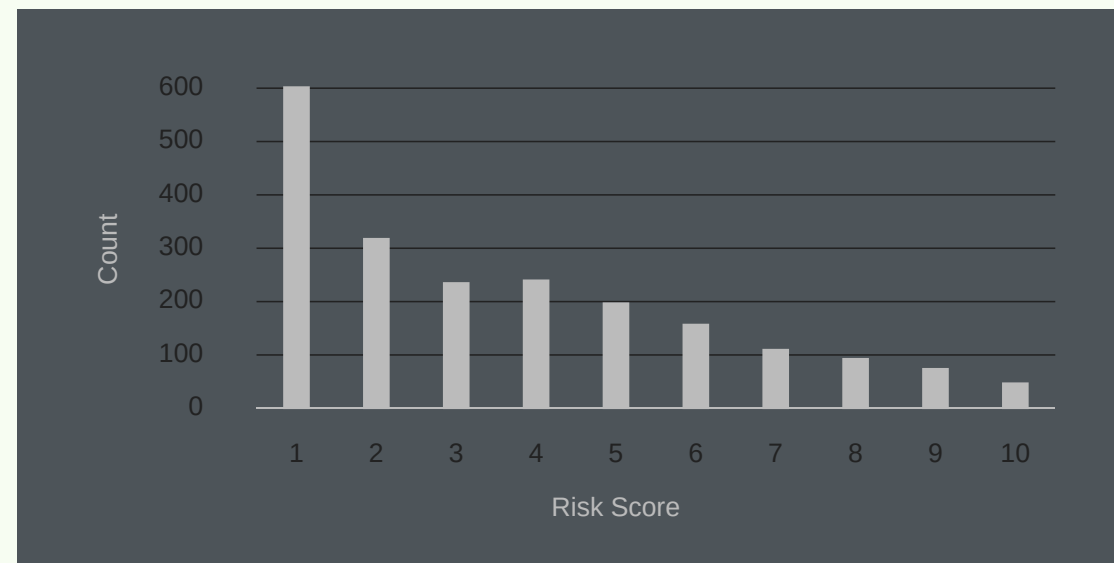
**Алън Тюринг, "Изчислителни машини и интелигентност"**

# Propublica <-> COMPASS

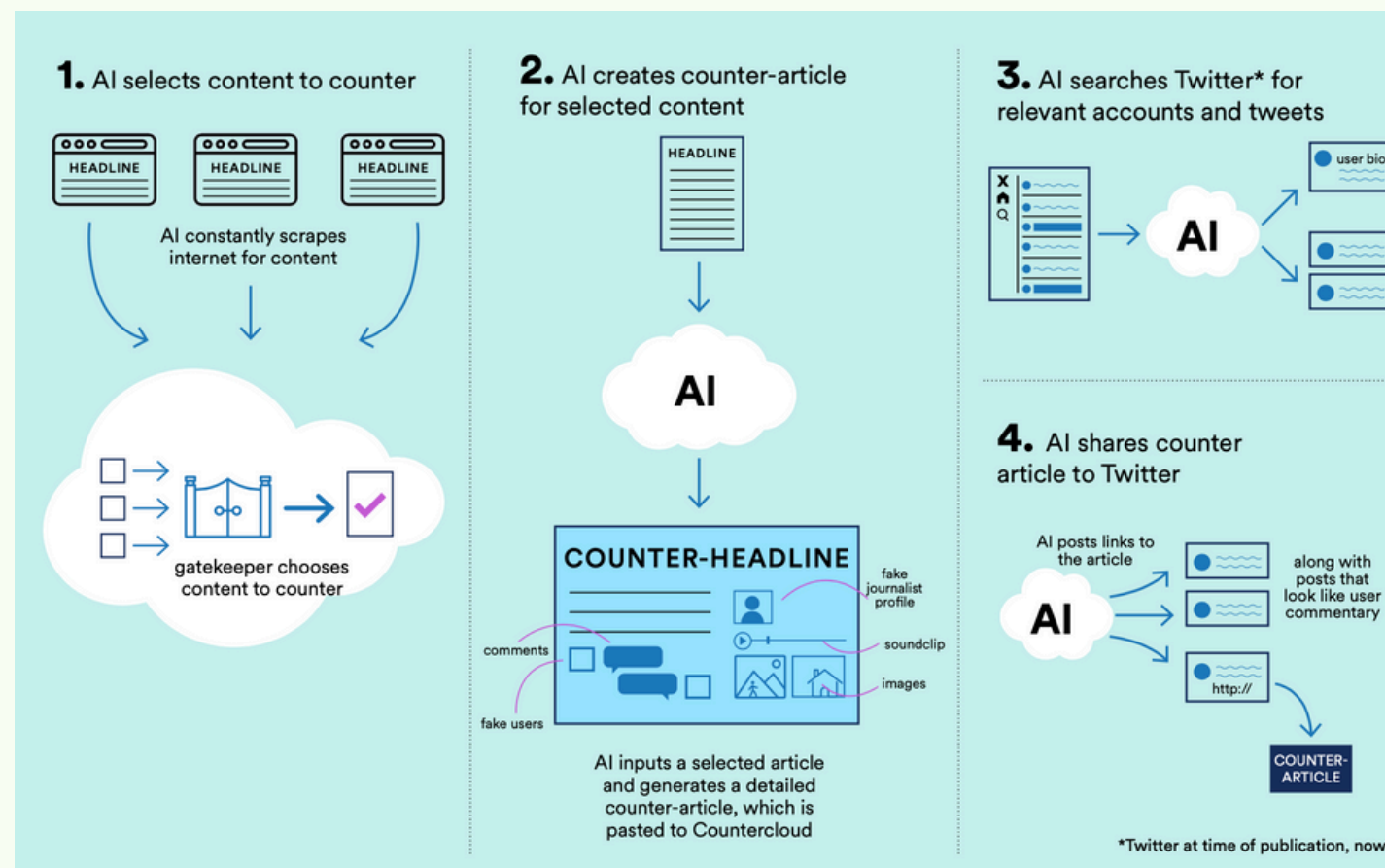
Концепцията за справедливост в изкуствения интелект включва равнопоставено отношение към индивидите от различни демографски групи.

Справедливостта е от съществено значение в **области с висок риск като правосъдие, медицина, изборен процес, заетост, енергетика, финанси.**

Въпреки усилията, много системи базирани на изкуствен интелект показват **предразсъдъци и водят до дискриминация.**



Сравнение между оценката на риска



Употреба в търсене на влияние в изборния процес

J. Angwin, S. Mattu, J. Larson, & L. Kirchner. (2016). Machine bias: There's software used across the country to predict future criminals. And it's biased against blacks. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

Nestor Maslej, et. al, "The AI Index 2024 Annual Report," AI Index Steering Committee, Institute for Human-Centered AI, Stanford University, Stanford, CA, April 2024

# Етични норми и тяхното влияние 1

Какво е реалното въздействие на етичните насоки върху човешкото вземане на решения в областта на изкуствения интелект и машинното обучение?

## European Commission's Ethics Guidelines for Trustworthy AI (2019)

Фокусира се върху осигуряване на законност, етика и устойчивост на ИИ системите.

## OECD Principles on AI (2019)

Подчертава необходимостта от растеж, устойчивото развитие и благосъстоянието.

## Beijing AI Principles (2019)

Призовава за отговорна употреба на ИИ за насърчаване на социалното добро.

## IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems (2019)

Предоставя цялостна рамка за етични съображения в ИИ.

## AI Now Report (2018)

Подчертава необходимостта от отговорност и прозрачност в ИИ системите.

## Partnership on AI (2018)

Сътрудничество между основни технологични компании за насърчаване на добри практики в ИИ.

## Montreal Declaration for Responsible AI (2017)

Цели да ръководи развитието на ИИ по начин, който уважава правата на човека.

## Asilomar AI Principles (2017)

Фокусира се върху дългосрочната безопасност и етичните съображения на ИИ.

## Future of Life Institute's AI Principles (2017)

Призовава за безопасна и полезна употреба на ИИ технологиите.

## UK House of Lords AI Report (2018)

Обсъжда последствията от ИИ за обществото и необходимостта от регулация.

## AI Ethics Guidelines by the Canadian Government (2019)

Подчертава прозрачността, отговорността и справедливостта в ИИ системите.

## Singapore's Model AI Governance Framework (2019)

Предоставя рамка за отговорно управление на ИИ.

## UNESCO's Recommendation on the Ethics of AI (2021)

Цели да насърчи глобалното етично развитие на ИИ.

## World Economic Forum AI Principles (2018)

Фокусира се върху социалното въздействие на ИИ и необходимостта от етични рамки.

## OpenAI Charter (2018)

Очертава принципи за отговорно развитие на ИИ технологиите.

## UAI4People's Ethical Framework (2018)

Предлага рамка за етичен ИИ, която приоритизира благосъстоянието на хората.

## The Montreal AI Ethics Institute's Guidelines (2019)

Цели да насърчи етичните практики в ИИ чрез ангажираност на общността.

## The Data Ethics Framework by the UK Government (2018)

Предоставя насоки за етично използване на данни в ИИ приложения.

## The AI Ethics Guidelines by the Australian Government (2020)

Фокусира се върху отговорната употреба на ИИ в публичната политика.

## The European Parliament's Resolution on AI (2020)

Призовава за цялостна регулаторна рамка за ИИ технологиите.

## The AI Ethics Guidelines by the German Government (2020)

Подчертава важността на човешкия надзор в ИИ системите.

## The AI Ethics Guidelines by the French Government (2020)

Призовава за етичен ИИ, който уважава основните права и свободи.

# Етични норми и тяхното влияние 2

Имат ли тези етични насоки реално въздействие върху човешкото вземане на решения в областта на изкуствения интелект и машинното обучение?

## Общи насоки в етичните норми

1. Отговорност
2. Прозрачност
3. Справедливост и равенство
4. Човекоцентричен дизайн
5. Поверителност и защита на данните
6. Безопасност и сигурност
7. Сътрудничество и ангажираност на заинтересованите страни
8. Непрекъснато наблюдение и оценка

Настоящите етични насоки рядко оказват значително влияние върху вземането на решения в AI индустрията. Често тези насоки служат като PR инструмент, без реално да променят практиките на компаниите.

## Основни аспекти

1. Ефективност на етичните насоки в AI.
2. Пропуски в съществуващите насоки.
3. Технически и етични предизвикателства.
4. Препоръки за подобряване на насоките.

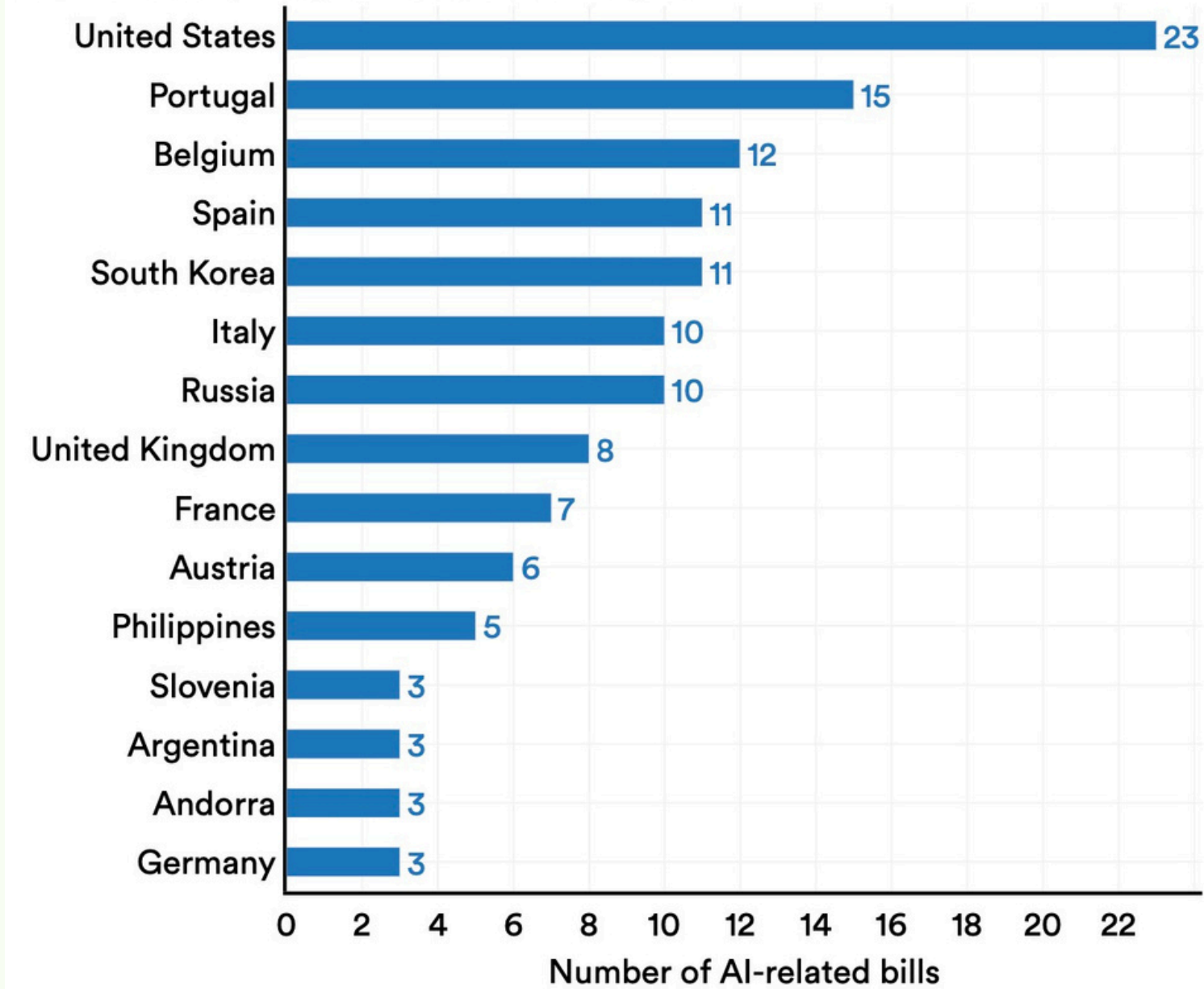
## Конфликтни точки

1. Точност <-> справедливост
2. Точност <-> обяснимост
3. Поверителност <-> прозрачност
4. Качество на услугите <-> поверителност
5. Персонализация <-> солидарност
6. Удобство <-> достойнство
7. Ефикасност <-> безопасност и устойчивост
8. Удовлетвореност <-> равенство

# Законодателство 2016 - 2023 географско покритие

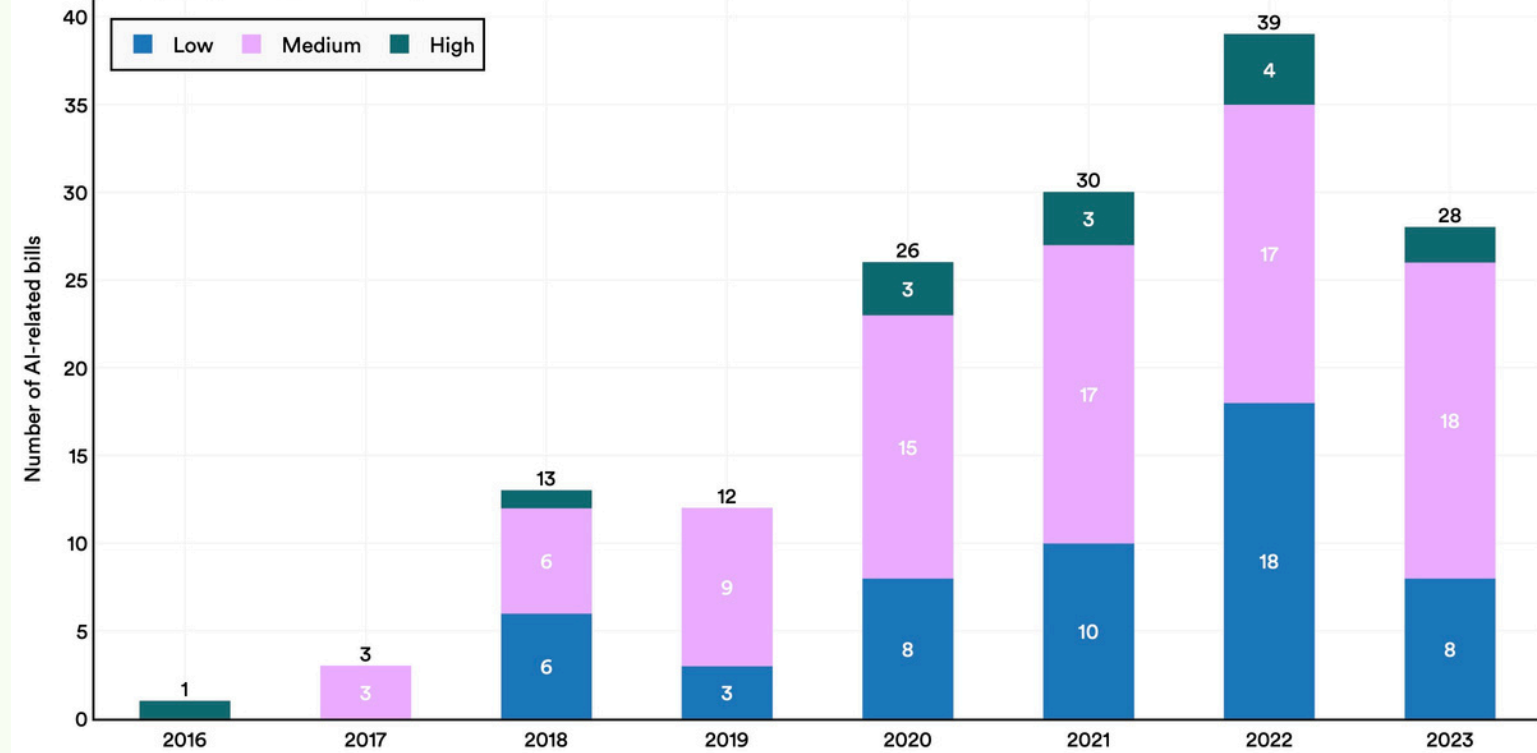
## Number of AI-related bills passed into law in select countries, 2016–23 (sum)

Source: AI Index, 2024 | Chart: 2024 AI Index report



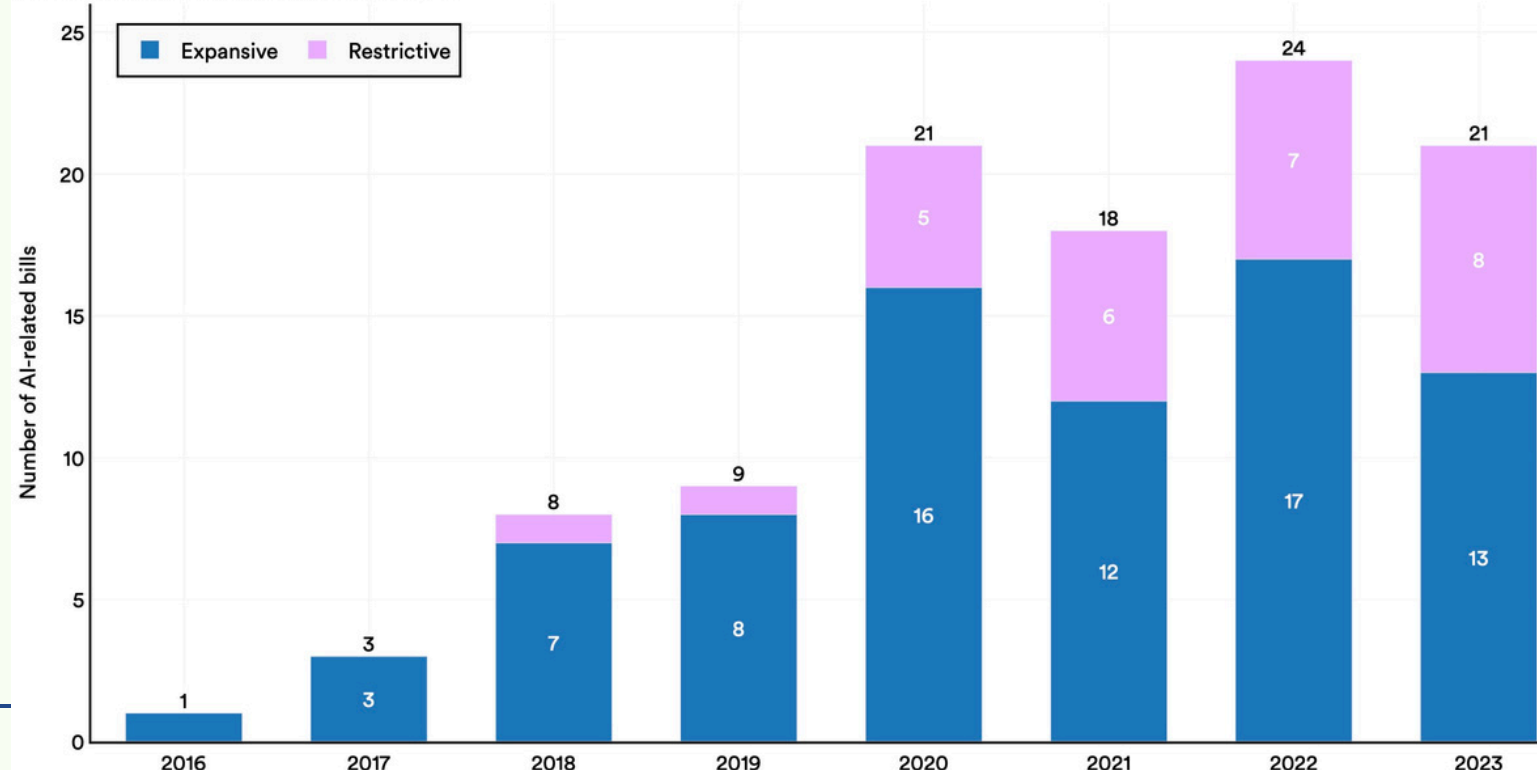
## Number of AI-related bills passed into law in select countries by relevance to AI, 2016–23

Source: AI Index, 2024 | Chart: 2024 AI Index report



## Number of AI-related bills passed into law in select countries by approach, 2016–23

Source: AI Index, 2024 | Chart: 2024 AI Index report

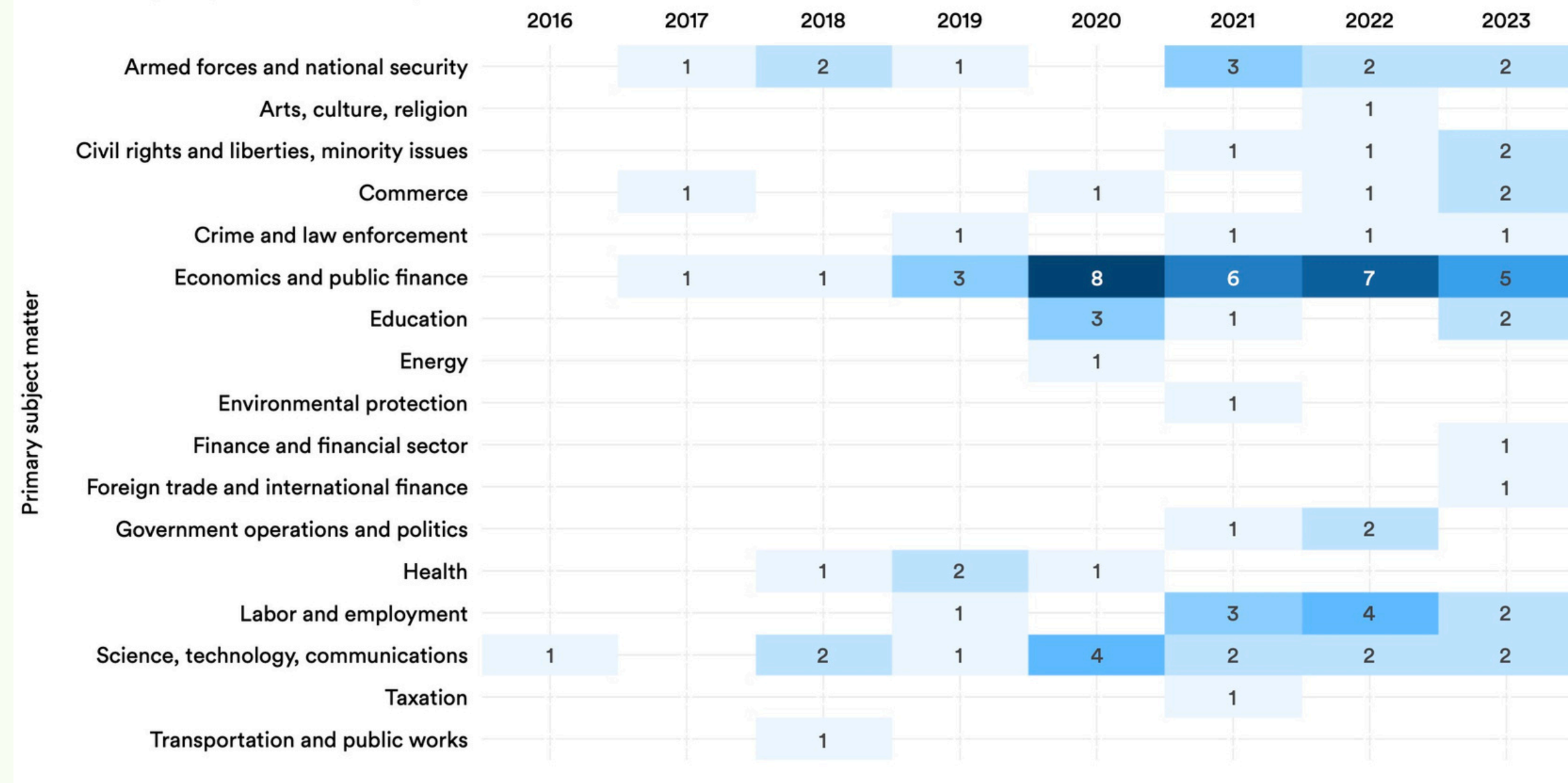




# Законодателство 2016 - 2023 по тематични области

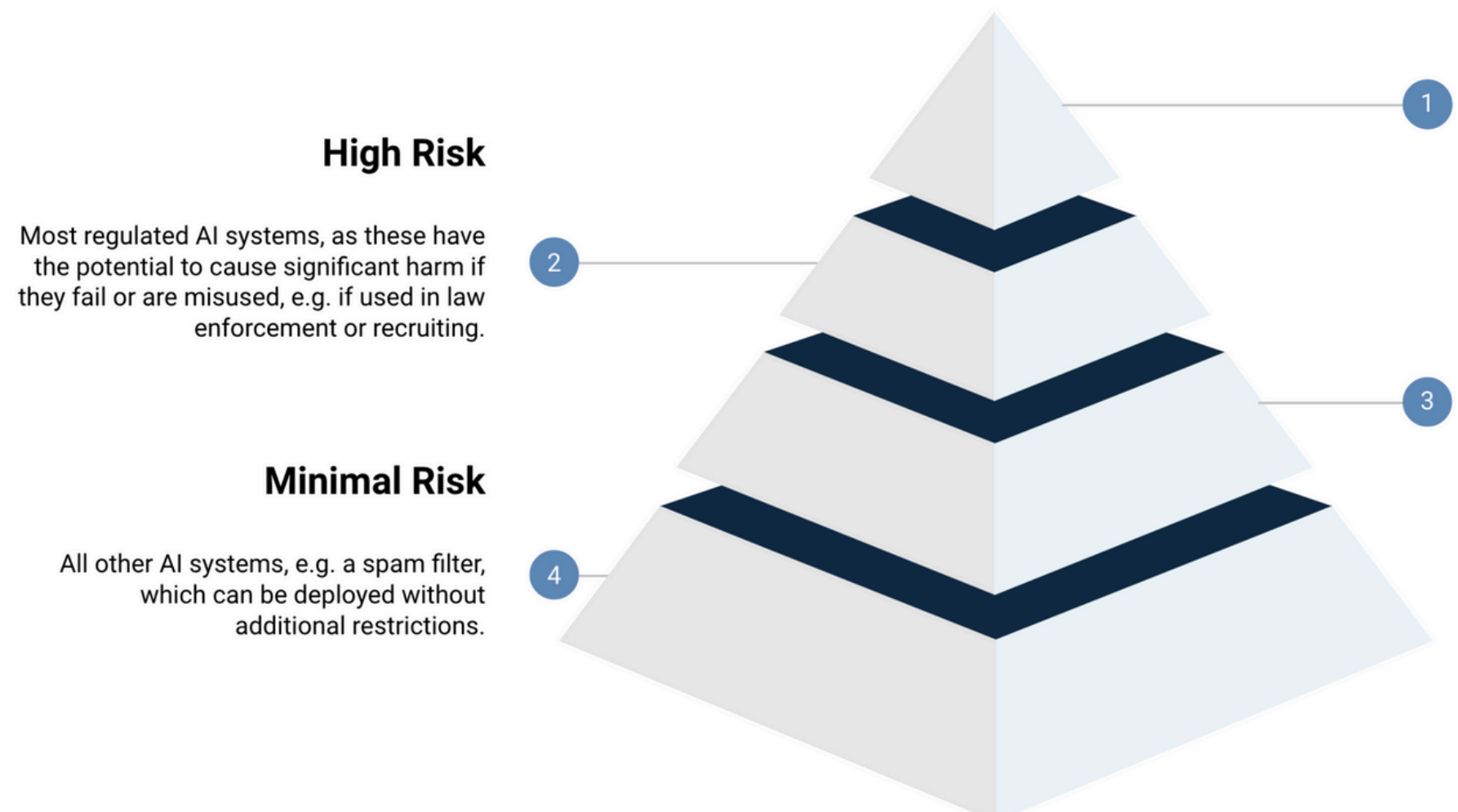
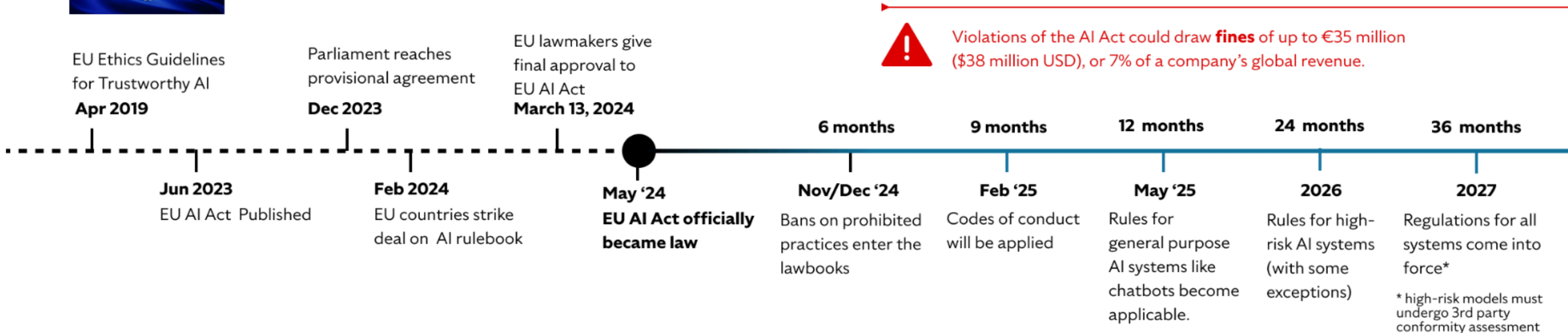
**Number of AI-related bills passed into law in select countries by primary subject matter, 2016–23**

Source: AI Index, 2024 | Chart: 2024 AI Index report





# Timeline of the EU AI Act



## Unacceptable Risk

1 Highest level of risk prohibited in the EU. Includes AI systems using e.g. subliminal manipulation or general social scoring.

## Limited Risk

3 Includes AI systems with a risk of manipulation or deceit, e.g. chatbots or emotion recognition systems. Humans must be informed about their interaction with the AI.



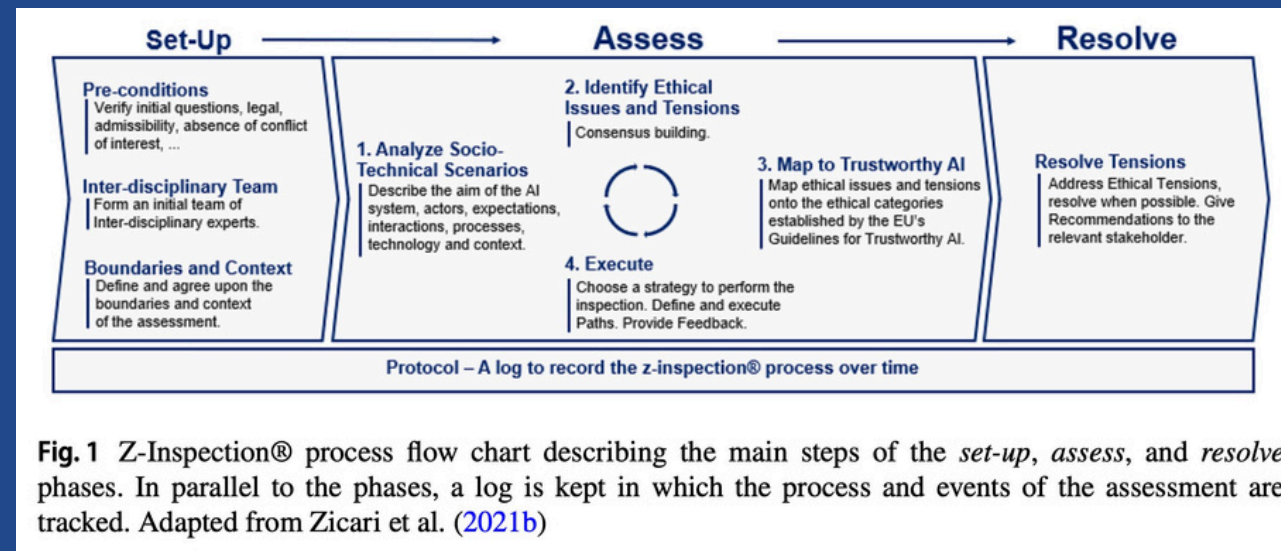
# Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self-assessment



## Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self-assessment

On the 17 of July 2020, the High-Level Expert Group on Artificial Intelligence (AI HLEG) presented their final Assessment List for Trustworthy Artificial Intelligence.

Shaping Europe's digital future



**Fig. 1** Z-Inspection® process flow chart describing the main steps of the *set-up*, *assess*, and *resolve* phases. In parallel to the phases, a log is kept in which the process and events of the assessment are tracked. Adapted from Zicari et al. (2021b)

### Table of Contents

|  |    |
|--|----|
| <b>Introduction</b>  | 3  |
| How to use this Assessment List for Trustworthy AI (ALTAI)       | 4  |
| <b>REQUIREMENT #1 Human Agency and Oversight</b>                 | 7  |
| Human Agency and Autonomy  | 7  |
| Human Oversight  | 8  |
| <b>REQUIREMENT #2 Technical Robustness and Safety</b>            | 9  |
| Resilience to Attack and Security                                | 9  |
| General Safety   | 9  |
| Accuracy   | 10 |
| Reliability, Fall-back plans and Reproducibility                 | 10 |
| <b>REQUIREMENT #3 Privacy and Data Governance</b>                | 12 |
| Privacy  | 12 |
| Data Governance  | 12 |
| <b>REQUIREMENT #4 Transparency</b>                               | 14 |
| Traceability   | 14 |
| Explainability   | 14 |
| Communication  | 15 |
| <b>REQUIREMENT #5 Diversity, Non-discrimination and Fairness</b> | 15 |
| Avoidance of Unfair Bias   | 16 |
| Accessibility and Universal Design                               | 17 |
| Stakeholder Participation  | 18 |
| <b>REQUIREMENT #6 Societal and Environmental Well-being</b>      | 18 |
| Environmental Well-being   | 19 |
| Impact on Work and Skills  | 19 |
| Impact on Society at large or Democracy                          | 20 |
| <b>REQUIREMENT #7 Accountability</b>                             | 21 |
| Auditability   | 21 |
| Risk Management  | 21 |
| <b>Glossary</b>  | 23 |

# Време за дискусия!