# APPENDIX B: TAXONOMY OF DATA COMMONS USE CASES FOR AI

*To supplement The Open Data Policy Lab's Blueprint to Unlock New Data Commons for AI, this document provides two frameworks that organizations that have or steward data can use to advance new data commons to enable public-interest AI use cases. This includes, for instance, knowledge institutions (such as libraries, universities, research centers) and others that hold unique and high-quality information assets (such as governments, non-profits).*

*The first framework maps the types of data that can be useful for generative AI applications. The second framework outlines 10 distinct AI use cases where a data commons approach could be valuable for society. These frameworks are a starting point only and aim to be exemplary, not all inclusive. We hope these resources inspire new data commons initiatives that help unlock AI's full potential.*

## CONTEXT

To organize the energy around data commons in a productive fashion, The Open Data Policy Lab conducted a series of structured analytic exercises to ideate potential use cases.

This effort began with a mapping exercise, identifying sectors where data commons for AI already exist and where there was substantial interest in securing more access to data. We relied on a search on the web (e.g. Google Scholar, news outlets, The Open Data Policy Lab's *Observatory of Generative AI and Open Data Use Cases*, and The GovLab's *Data Collaboratives Explorer*) and our network of experts (e.g. Open Data Action Labs). This review helped us understand current efforts across industries. We summarized this work within the blog: *The Emergent Landscape of Data Commons: A Brief Survey and Comparison of Existing Initiatives*.

We then conducted another mapping of data types that could be useful for different aspects of generative AI, drawing from Microsoft's *Government Data Commons: A Technical Guide*, submissions to the *United States Department of Commerce Request for Information on AI-Ready Open Government Data Assets*, The Data Provenance Initiative's *A Large Scale Audit of Dataset Licensing & Attribution in AI*, the Open Data Institute's *Data for AI Taxonomy*, *Papers with Code*, and our network of experts.

We focused on data types that were cited as valuable for generative AI or where authors expressed there is a need for data. We chose not to include data types with clear ethical and data responsibility issues (data types such as electronic health records that have been widely criticized by the public for its use within AI systems). We also tried to avoid data types where extensive data commons already exist to avoid the duplication of efforts. Using this information, we developed a taxonomy to organize these data types.

At the same time, we conducted a "What if?" Analysis in which we asked ourselves to consider the types of use cases where new AI data commons could benefit society. Drawing on the information gathered thus far, we synthesized 10 categories of use cases for new data commons in the context of AI.
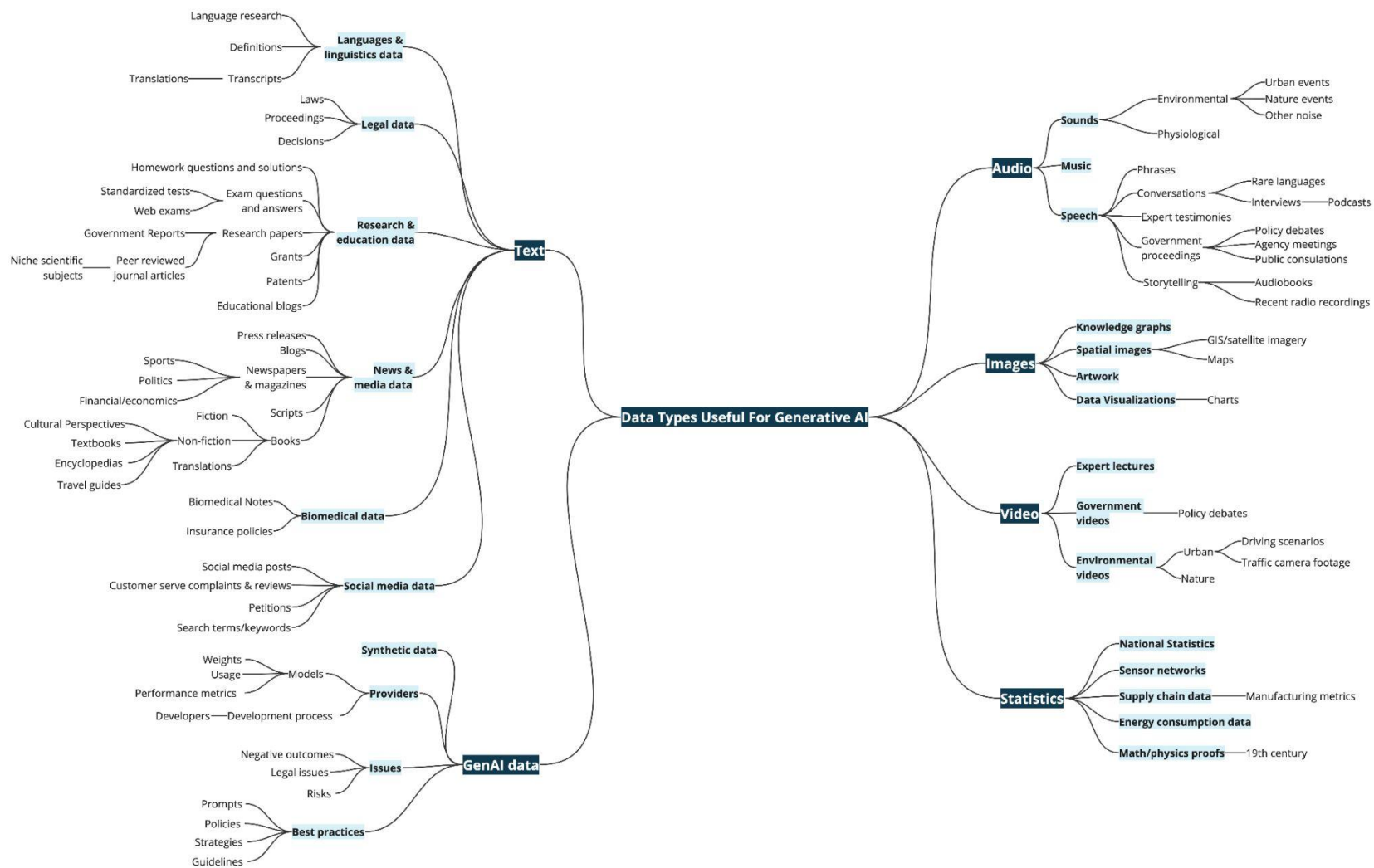
In what follows we provide a summary of these exercises. While these exercises primarily focus on generative AI, it is important to note that there is still a need for data commons for machine learning and neutral networks as well.

## TAXONOMY OF DATA TYPES

Below we provide a list of data types that could be valuable for data commons for AI. These data types are divided into six main categories. Creating data commons that are truly fair and beneficial to all requires sourcing data responsibly. That being said, the data commons should only include data sources with appropriate licensing.

| Category | Data Types (Exemplary) |
|---|---|
| 1. Text | • Language and linguistics data: Language research, words, definitions, translations, and transcripts<br>• Legal data: Laws, legal proceedings, and court case decisions<br>• Research and education data: Exam questions and answers, educational blogs, homework questions and solutions, research reports, peer-reviewed journal articles (focusing on niche scientific subjects), grants, and patents<br>• News and media data: Press releases, blogs, newspapers and magazines, scripts, books (e.g. textbooks, cultural perspectives, encyclopedias, travel guides, and translations)<br>• Biomedical data: Notes and insurance policies<br>• Social media data: Social media posts, customer reviews, online petitions, search terms/keywords |
| 2. Audio | • Sounds: Environmental sounds (urban and nature) and physiological sounds (e.g. breathing patterns)<br>• Music<br>• Speech Data: Conversations, expert testimonies, phrases, government proceedings (e.g. policy debates, consultations), storytelling data (e.g. audiobooks, recent radio recordings) |
| 3. Images | • Knowledge graphs<br>• Spatial images: GIS/satellite data, maps<br>• Artwork (including photographs)<br>• Data visualizations: Charts, tables, infographics |
| 4. Video | • Expert lectures<br>• Government videos: Policy debates<br>• Environmental videos: Urban videos (e.g. driving scenarios, traffic camera footage) and nature videos |

| 5. Statistics | <ul><li>National statistics</li><li>Sensor networks</li><li>Supply chain data: Manufacturing metrics</li><li>Energy consumption data</li><li>Math/physics proofs: 19th century physics proofs</li></ul> |
|---|---|
| 6. Generative AI data | <ul><li>Synthetically created datasets</li><li>Data about generative AI providers: Models (weights, usage, performance), development processes</li><li>Issues associated with generative AI: Unwanted outcomes, legal cases, other risks</li><li>Best practices when using generative AI: Prompts, policies, strategies, and guidelines</li></ul> |

**Figure 2. Summary of Data Types**
*View the full map HERE*

## TAXONOMY OF USE CASES

The following chart summarizes our taxonomy of data commons use cases that could benefit society. In this taxonomy, we focus on the purpose for developing AI technologies, the types of technologies that could be developed and the scope of the use or who it is affecting most. We also outline potential stakeholders that may be involved. This list of stakeholders is exemplary only and non-exhaustive.

| Area of Focus | Purpose | Scope | Potential Outcomes | Stakeholders |
|---|---|---|---|---|
| 1. **Research and Scientific Discovery** | Enhance scientific exploration and accelerate new discoveries and innovations | Societal | Drug discovery, hypothesis generation, climate modeling, open science collaboration | • Scientists and researchers<br>• Research and academic institutions<br>• Research funders<br>• Publishers/journals<br>• Government research agencies<br>• Multilateral organizations<br>• Representatives of the public |
| 2. **Education, Training, and Learning Support** | Support education and learning through personalized and interactive tools | Individual | Personalized learning platforms, training simulations, assistive tools, content generation | • Students<br>• Educators<br>• Research and academic institutions<br>• Education technology companies<br>• Government education departments<br>• Multilateral organizations |

| 3. Content Development and Knowledge Preservation | Create and tailor content for education, service delivery, communications, and the preservation of knowledge across cultures | Organizational | Multilingual content, public health campaigns, digital heritage, AI-assisted infographics/visual communications | • Archives, museums, and libraries<br>• Government communications departments<br>• Civil society organizations (NGOs, advocacy groups, and unions)<br>• Multilateral organizations<br>• Representatives of the public |
|---|---|---|---|---|
| 4. Modeling, Simulation, and Anticipatory Decision-Making | Anticipate future outcomes through realistic, data-driven scenario generation | Societal | Climate simulations, policy modeling, urban planning (e.g. digital twins of cities), crisis management, multi-stakeholder decision models, risk analysis | • Foresight and forecasting practitioners<br>• Government strategists and policy makers<br>• Local governments<br>• Philanthropies<br>• Climate agencies<br>• Representatives of the public |
| 5. Pattern Recognition, Inference, and Insight Generation | Analyze and identify trends and anomalies in data for better decision-making | Societal | Traffic optimization, crime pattern simulation, disinformation detection, economic trends | • Research and academic institutions<br>• Government agencies<br>• National Statistics Offices<br>• Local governments<br>• Multilateral organizations<br>• Civil society organizations |

| | | | | |
|---|---|---|---|---|
| 6. **Personalization and Customization of Services** | Tailor services and experiences to individual needs | Individual | Personalized healthcare, adaptive public services, civic engagement tools | • Domain experts<br>• Government agencies<br>• DEI departments<br>• Civic tech<br>• NGOs and multilateral organizations |
| 7. **Automation of Repetitive Tasks** | Increase efficiency by automating time-consuming processes | Organizational | Automated reporting, legal document drafting, customer support chatbots | • Government agencies<br>• Multilateral organizations<br>• NGOs and non-profits<br>• Philanthropies<br>• Advocacy groups<br>• Businesses |
| 8. **Synthetic Data Generation** | Preserve privacy while creating synthetic data for model training and research | Organizational | Synthetic datasets for healthcare, research benchmarks, addressing data scarcity | • PII departments<br>• National Statistics Offices<br>• Multilateral organizations |
| 9. **Real-Time Insights and Adaptive Systems** | Provide feedback in real-time for dynamic system improvement | Societal | Real-time traffic and health monitoring, adaptive disaster response | • Government agencies<br>• Multilateral organizations<br>• Emergency response departments<br>• Local governments |

| | | | | |
|---|---|---|---|---|
| **10. Collaborative Intelligence** | Enhance collective problem-solving and finding consensus among diverse groups | Organizational | Mediation and conflict resolution, administrative support, public consultation tools | • Government agencies<br>• Multilateral organizations and NGOs<br>• Policy makers<br>• Civil society organizations<br>• Representatives of the public |