



HAL
open science

Towards a multi-criteria evaluation of the environmental footprint of generative ai services

Adrien Berthelot, Eddy Caron, Mathilde Jay, Laurent Lefèvre

► To cite this version:

Adrien Berthelot, Eddy Caron, Mathilde Jay, Laurent Lefèvre. Towards a multi-criteria evaluation of the environmental footprint of generative ai services. ICT4S 2024 - International Conference on Information and Communications Technology for Sustainability, Jun 2024, Stockholm, Sweden. pp.1-1, 2024. hal-04586653

HAL Id: hal-04586653

<https://inria.hal.science/hal-04586653v1>

Submitted on 24 May 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - ShareAlike 4.0 International License

TOWARDS A MULTI-CRITERIA EVALUATION OF THE ENVIRONMENTAL FOOTPRINT OF GENERATIVE AI SERVICES

Adrien Berthelot^{1,2}, Eddy Caron¹, Mathilde Jay^{1,3} and Laurent Lefèvre¹

¹Univ. Lyon 1 UCBL, EnsL, CNRS, Inria, LIP. Lyon, France

²OCTO Technology. Paris, France

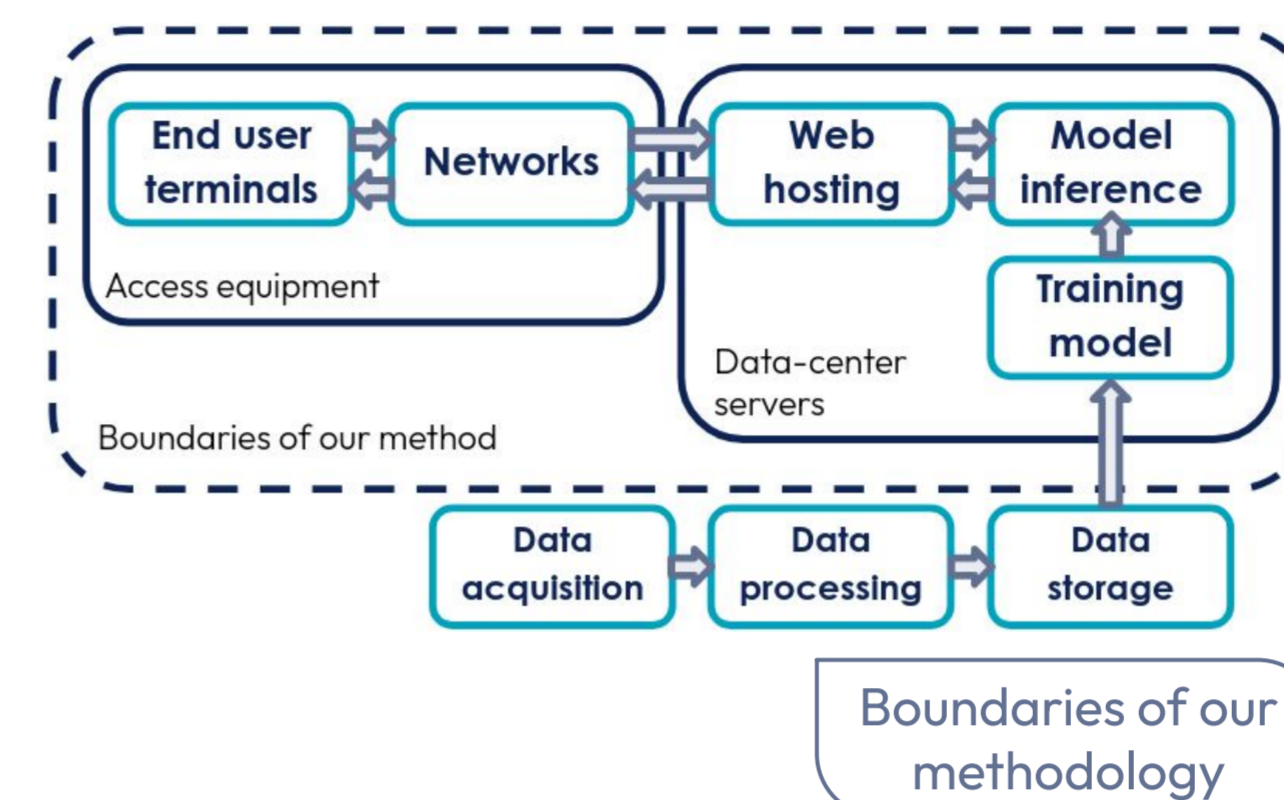
³Univ. Grenoble Alpes, Grenoble INP, LIG. Grenoble, France



Introduction

Generative AI represents a new stage in digital transformation through its many applications. Unfortunately, by accelerating the growth of digital technology, Gen-AI is contributing to the **multiple environmental damages** caused by its sector.

We propose to study not only the impact of developing a model but also that of its deployment and use as a **service**. The figure at the left shows what we consider the standard structure of a Gen-AI service. Our methodology (presented in the diagram below) is applied and validated on an AI service based on **Stable Diffusion**, an open-source text-to-image generative deep-learning model.

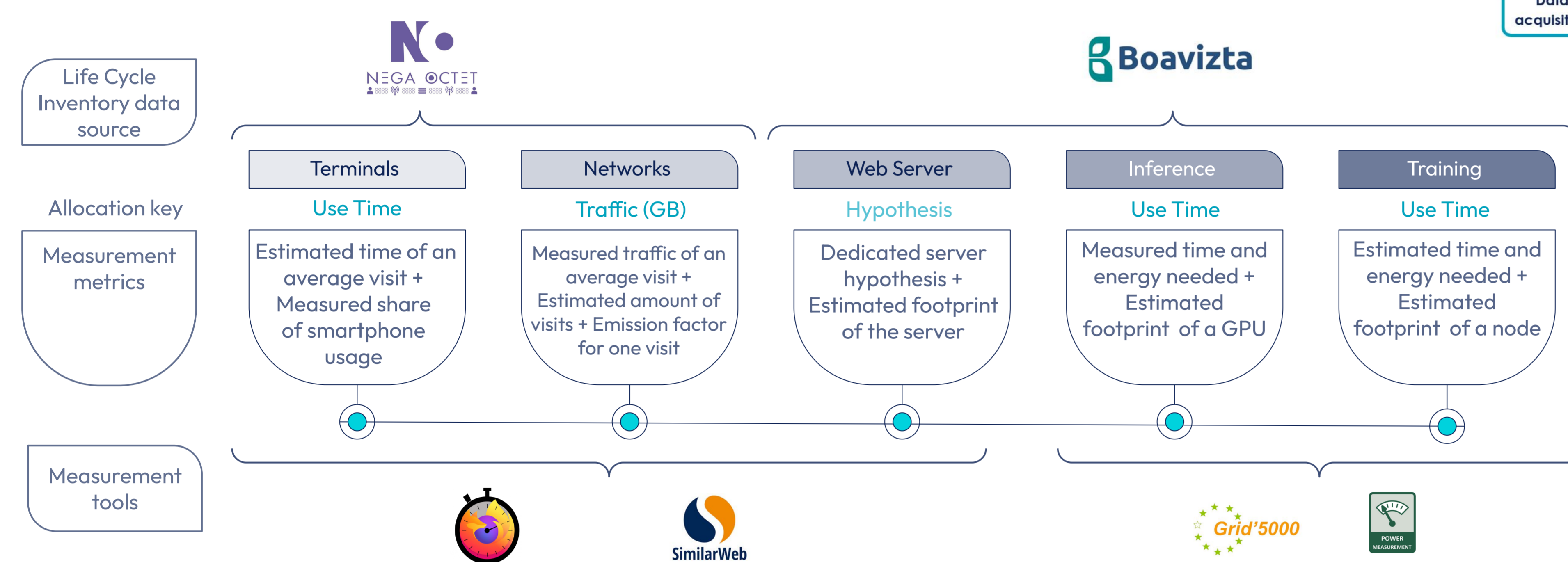


Training electricity cost

Replication is the most reliable method to estimate the electricity consumption. But for training, it would be too expensive.

Our solution: replicate a fraction of it and use **linear regression to extrapolate** the results.

The figure below shows the difference between the methods used to estimate the electricity consumption. It leads to a **better understanding of electricity consumption**.

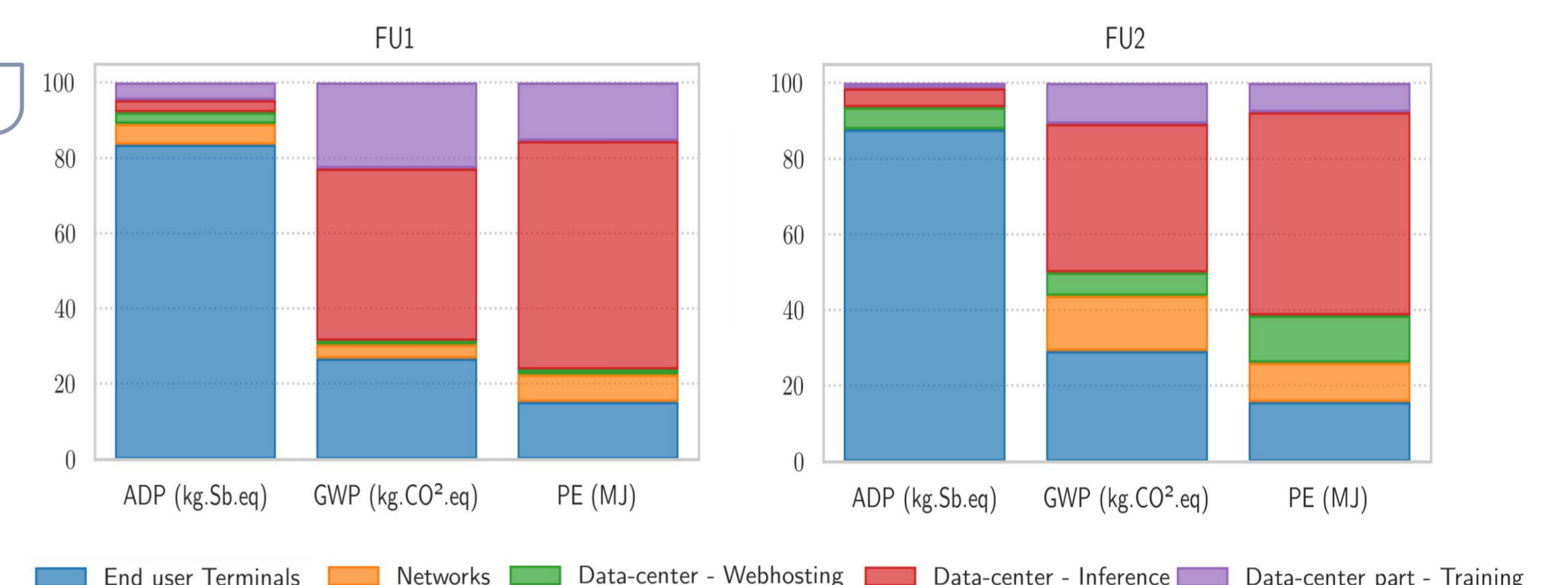


Environmental footprint of the service

Functional Unit (FU)	Abiotic Depletion Potential (kg.Sb.eq)	Warming Potential (kg.CO2.eq)	Primary Energy (MJ)
FU1 - Average single use of service	6.72e ⁻⁸	7.84e ⁻³	2.02e ⁻¹
FU2 - A year of service	4.64	3.60e ⁺⁵	8.93e ⁺⁶

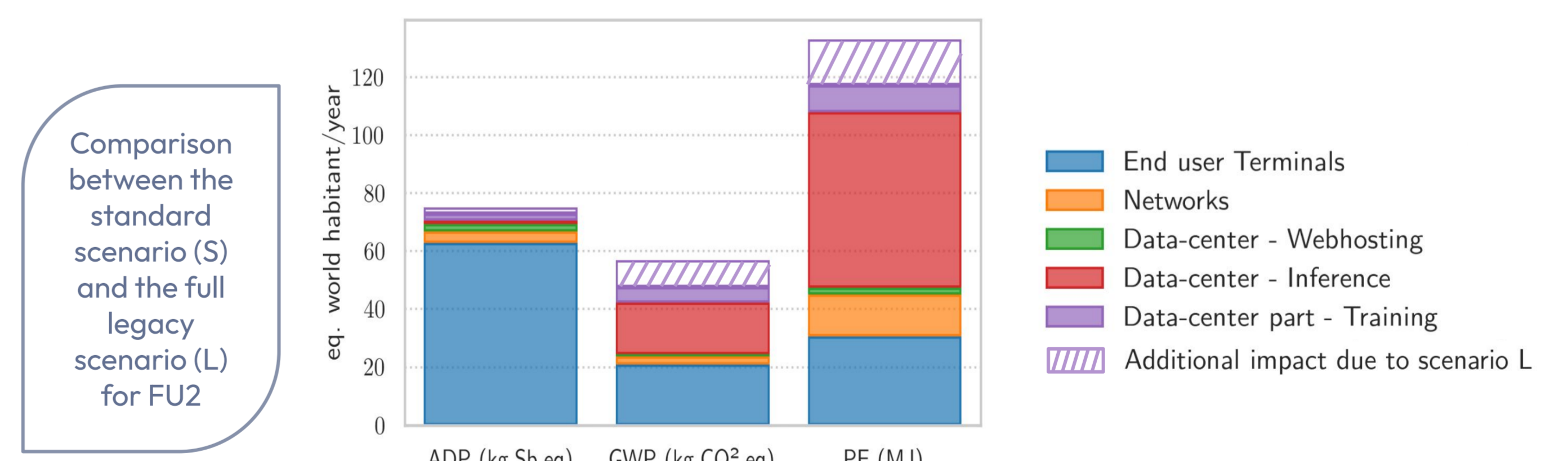
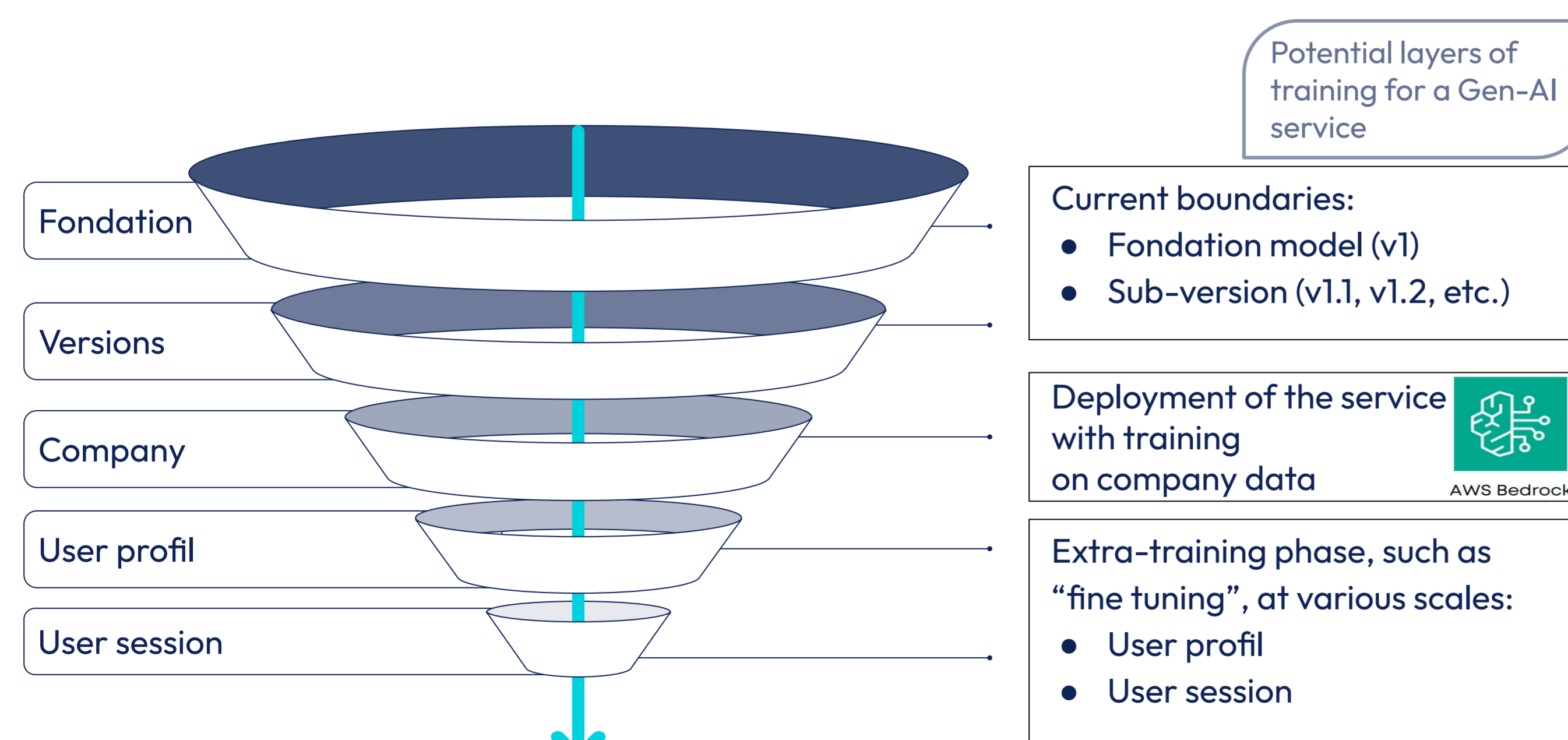
Impact distributions for FU1 and FU2

Environmental impact of Stable Diffusion for FU1 and FU2



Sensitivity to the utilization rate and the scope of training

- Average utilization rate (**AUR**) of data center equipment
 - AUR = percentage of time in the equipment lifespan during which it is actively used, as opposed to when it is either idle or on standby.
 - Crucial to the allocation for cost-heavy data center equipment.
 - Impact remains small if the AUR is higher than **20%** (is that reasonable?)
- Scope of training phases taken into account.
 - Diagram below shows hierarchical **fine-tuning layers** in training. Our standard scenario (S) only includes versions v1.1 and v1.2.
 - We consider versions up to v1.5 in a full legacy scenario (L). Figure on the right shows the additional impact due to these additional versions.



Conclusions

- Generative AI offers digital services that are particularly costly in environmental terms.
- The total footprint is not concentrated in a single part and a single impact.
- A large proportion of GHG emissions can be avoided, but this will **not be enough**.
- The transformation of data centers induced by the multiplication of these services will generate numerous impacts of 2^e and 3^e order.
- More than generative AI as a technology, it is the rapid, growing, uncontrolled **deployment** of it as a service that represents a problem for our environment.

