

Targeted consultation in preparation of the COMMISSION GUIDELINES TO CLARIFY THE SCOPE OF THE OBLIGATIONS OF PROVIDERS OF GENERAL-PURPOSE AI MODELS IN THE AI ACT

Disclaimer: This is a working document by the AI Office for consultation and does not prejudice the final decision that the Commission may take on the guidelines. The responses to this consultation will provide important input to the Commission when preparing the guidelines.

1. BACKGROUND AND PURPOSE OF THE CONSULTATION

The AI Office is launching a multi-stakeholder consultation to assist in the preparation of guidelines on general-purpose AI which aim at clarifying the scope of the rules for providers of general-purpose AI models in Regulation (EU) 2024/1689 ('AI Act'). Those rules will enter into application on 2 August 2025.

Artificial Intelligence ("AI") promises huge [benefits](#) to our economy and society. General-purpose AI models play an important role in that regard, as they can be used for a variety of tasks and may form the basis for a range of downstream AI systems. The [AI Act](#) aims to ensure that general-purpose AI models are transparent, safe, and trustworthy.

The AI Office is dedicated to facilitating compliance of providers of general-purpose AI models with their obligations under the AI Act. To this end, the Commission guidelines on general-purpose AI are expected to clarify key concepts in the AI Act, such as what is a 'general-purpose AI model', a 'provider of a general-purpose AI model', a 'placing on the market of a general-purpose AI model', and how to estimate the computational resources used for training a general-purpose AI model. Beyond these conceptual clarifications, the guidelines are expected to clarify how the AI Office will work with providers who must comply with the general-purpose AI rules, to support them in their compliance. The Commission's Joint Research Centre is providing scientific evidence to these guidelines.

The Commission guidelines on general-purpose AI will complement the General-Purpose AI Code of Practice ('Code') which will set out commitments to which providers of general-purpose AI models may adhere to ensure compliance with their obligations under the AI Act. More specifically, the Code will detail how those providers can ensure compliance with the documentation and copyright obligations that apply to all providers of general-purpose AI models in Article 53 AI Act. For the most advanced models that may pose systemic risk, the Code will outline how providers of those models can ensure compliance with the systemic risk assessment and mitigation obligations throughout the lifecycle of the model in Article 55 AI Act. Building on an initial multi-stakeholder consultation and several feedback rounds, the Code is currently being drafted by independent experts, with input from over 1,000 stakeholders across industry, civil society, academia, and others. The Code will align with emerging industry best practices and international approaches, ensuring that providers can innovate, while maintaining the trust of consumers and SMEs in the technology they use.

Both the Commission guidelines on general-purpose AI and the final General-Purpose AI Code of Practice are expected to be published in May or June 2025. In parallel, the Commission will also publish separate guidelines containing the template for the summary of the content used for training, to assist providers of general-purpose AI models to fulfil their

obligation in Article 53(1), point (d), AI Act. The Commission has already published guidelines on the [AI system definition](#) and the [prohibited AI practices](#) under the AI Act.

Although the Commission guidelines on general-purpose AI will be non-binding, they will provide important clarifications on how the Commission, which is exclusively responsible for the supervision and enforcement of the obligations of providers of general-purpose AI models under the AI Act, will interpret and apply those obligations when exercising its tasks under the AI Act. The guidelines are expected to evolve over time and will be updated as necessary, in particular in light of evolving technological developments. An authoritative interpretation of the AI Act may ultimately only be given by the Court of Justice of the European Union (CJEU).

This targeted consultation aims to gather a broad range of input and perspectives. We invite submissions from all stakeholders with relevant expertise and perspectives **through [this survey](#)¹ by Thursday, 22 May 2025, 12:00 (noon) CET**, particularly from industry actors such as providers of general-purpose AI models and downstream providers of AI systems built on those models, civil society, academia, other independent experts, and public authorities.

2. PRELIMINARY OVERVIEW OF THE CONTENT OF THE GUIDELINES

These Commission guidelines on general-purpose AI are expected to cover the following topics:

- What is a general-purpose AI model (Section 3.1)
- Who is the provider of a general-purpose AI model, and when is a downstream modifier a provider (Section 3.2)
- What constitutes a placing on the market of a general-purpose AI model, and when do the open-source exemptions apply (Section 3.3)
- Estimating the computational resources used to train or modify a model (Section 3.4)
- Transitional rules, grandfathering, and retroactive compliance (Section 3.5)
- Effects of adherence to and signature of a code of practice (Section 3.6)
- Supervision and enforcement of the general-purpose AI rules (Section 3.7)

3. PRELIMINARY APPROACH FOR THE CONTENT OF THE GUIDELINES

3.1. General-purpose AI model

The definition of “general-purpose AI model” is key to understanding whether an entity must comply with the AI Act’s rules for general-purpose AI models.

Article 3(63) AI Act defines a “general-purpose AI model” as “*an AI model, including where such an AI model is trained with a large amount of data using self-supervision at scale, **that displays significant generality and is capable of competently performing a wide range of distinct tasks** regardless of the way the model is placed on the market and that can be integrated into a variety of downstream systems or applications, except AI models that are used for research, development or prototyping activities before they are placed on the market.*” *[emphasis added]*

¹ https://ec.europa.eu/eusurvey/runner/GPAI_Guidelines_Consultation_2025

The Commission guidelines will build on the guidance already provided by the following recitals:

- Recital 97 AI Act: *“The notion of general-purpose AI models should be clearly defined and set apart from the notion of AI systems to enable legal certainty. The definition should be based on the key functional characteristics of a general-purpose AI model, in particular the generality and the capability to competently perform a wide range of distinct tasks. These models are typically trained on large amounts of data, through various methods, such as self-supervised, unsupervised or reinforcement learning. (...) Although AI models are essential components of AI systems, they do not constitute AI systems on their own. AI models require the addition of further components, such as for example a user interface, to become AI systems. AI models are typically integrated into and form part of AI systems. (...) The definition should not cover AI models used before their placing on the market for the sole purpose of research, development and prototyping activities. This is without prejudice to the obligation to comply with this Regulation when, following such activities, a model is placed on the market.”*
- Recital 98 AI Act: *“Whereas the generality of a model could, inter alia, also be determined by a number of parameters, models with at least a billion of parameters and trained with a large amount of data using self-supervision at scale should be considered to display significant generality and to competently perform a wide range of distinctive tasks.”*
- Recital 99 AI Act: *“Large generative AI models are a typical example for a general-purpose AI model, given that they allow for flexible generation of content, such as in the form of text, audio, images or video, that can readily accommodate a wide range of distinctive tasks.”*

This section of the guidelines is expected to provide guidance on when a model should be considered a general-purpose AI model (Section 3.1.1), as well as clarify the difference between models and model versions (Section 3.1.2).

3.1.1 Conditions for sufficient generality and capabilities

The Commission guidelines are expected to outline practical conditions for determining when a model displays sufficient generality and capabilities to classify it as a general-purpose AI model. This part of the guidelines, in particular the training compute threshold, is likely to change in the future in light of evolving technological developments.

To determine whether a model is a general-purpose AI model, the decisive question is whether the model *“displays significant generality and is capable of competently performing a wide range of distinct tasks”* (Article 3(63) AI Act).

Given the wide variety of potential uses for general-purpose AI models, it is not possible to provide a precise list of tasks that determine whether a model *“displays significant generality and is capable of competently performing a wide range of distinct tasks”*. Whilst benchmarks and other tools to evaluate model capabilities and generality may be used in some cases to assess whether a model should be considered a general-purpose AI model, they are still too immature to form the basis of reliable criteria for classifying a model as a general-purpose AI model. Furthermore, a criterion that can easily be checked by potential providers is desirable in order to limit the burden on the many entities that will have to assess whether they are providers of a general-purpose AI model.

Recital 98 AI Act specifies that “*models with at least a billion of parameters and trained with a large amount of data using self-supervision at scale should be considered to display significant generality and to competently perform a wide range of distinctive tasks*”. However, recital 98 AI Act leaves unspecified what “a large amount of data” means. To capture model generality and capabilities, the AI Office’s preliminary approach is to set a threshold in terms of the amount of computational resources used to train a model (training compute). This approach combines number of parameters and amount of training data in a single number (training compute) which is roughly proportional to the product of its number of parameters² and its number of training examples,³ instead of measuring model size and training data size separately.

Training compute is an imperfect proxy for generality and capabilities but currently the best metric for legal certainty. The AI Office will continue to analyse availability of alternative solution(s) that could be used to assess generality and capabilities with relative ease, especially for smaller entities.

In particular, the AI Office’s preliminary approach is to presume that a model that can generate text and/or images is a general-purpose AI model if its training compute is greater than 10^{22} FLOP. This threshold is based on the fact that models with one billion parameters that can generate text and/or images are typically trained using approximately 10^{22} FLOP (see Annex A.1. for relevant examples and calculations). For general guidance on how training compute may be estimated, see Section 3.4.1.

Models that cannot generate text and/or images may be considered general-purpose AI models if they have a level of generality and capabilities comparable to such models.

The presumption that a model that can generate text and/or images is a general-purpose AI model based on its training compute is rebuttable. This means that if its training compute meets the threshold, the model is expected to have sufficient generality and capabilities to qualify as a general-purpose AI model unless there is evidence to the contrary. Whether a model displays significant generality and is capable of competently performing a wide range of distinct tasks depends not only on training compute, but also on the modality and other characteristics of the data used for training. For example, a model that is only usable for transcription of speech should not be considered capable of competently performing a wide range of tasks, even if its training compute meets the threshold. Thus, the modality or other specific characteristics of the training data may offer grounds for rebuttal of the presumption.

If the training compute of a model capable of generating text and/or images is lower than the threshold of 10^{22} FLOP, then the model is presumed to lack sufficient generality and capabilities to be a general-purpose AI model, unless there is evidence to the contrary.

Examples of models in scope:

- A model is trained on a broad range of natural language data curated and scraped from the internet (as is currently typical for language models) using 10^{23} FLOP.
 - The model is presumed to be a general-purpose AI model, because the model can generate text and its training compute meets the FLOP-threshold. Training

² More precisely, number of parameters active in each forward pass.

³ Where the training data is text, a training example corresponds to a training token.

on a broad range of natural language also indicates that the model should display significant generality and be capable of competently performing a wide range of distinct tasks, so it is unlikely that the presumption can be rebutted.

- A model is trained on speech data (i.e. audio) but no other modality, using 10^{23} FLOP. Benchmarks show that it is able to perform a similar range of tasks with a comparable level of performance as general-purpose AI models trained for text generation.
 - The model is not presumed to be a general-purpose AI model, because it is not trained to generate either text or images. However, the benchmarks indicate that it should be considered a general-purpose AI model.

Examples of models out of scope:

- A model is trained on natural language data, using 10^{20} FLOP.
 - The model is not presumed to be a general-purpose AI model, because its training compute is below the FLOP-threshold.
- A model is trained on a specialised data set containing only code but no other kinds of text, using 10^{23} FLOP.
 - The model is presumed to be a general-purpose AI model, because the model can generate text and its training compute meets the FLOP-threshold. However, this presumption may be rebutted since the model can only competently perform a narrow set of tasks (coding).
- A model is trained specifically for the task of transcribing speech to text, using 10^{23} FLOP.
 - The model is presumed to be a general-purpose AI model, because the model can generate text and its training compute meets the FLOP-threshold. However, this presumption may be rebutted since the model can only competently perform a narrow set of tasks (transcribing speech).
- A model is trained specifically for the task of increasing the resolution of images given as input, but not to generate images given a text description, using 10^{23} FLOP.
 - The model is presumed to be a general-purpose AI model, because the model can generate images and its training compute meets the FLOP-threshold. However, this presumption may be rebutted as the model can only competently perform a narrow set of tasks (upscaling images).

3.1.2 Differentiation between distinct models and model versions

The Commission guidelines are expected to provide details about when the provider of a general-purpose AI model is considered a ‘distinct model’ versus a version of the same model (‘model version’). This part of the guidelines, in particular the focus on the large pre-training run, is especially likely to change in the future in light of evolving technological developments.

The AI Office’s preliminary approach is to consider the large pre-training run the beginning of the lifecycle of a general-purpose AI model. In this context, a large pre-training run is understood as the foundational training run conducted on a large amount of data to build the model’s general capabilities, which may take place after smaller experimental training runs, and which may be followed by fine-tuning for specialization or other post-training enhancements. Therefore, the AI Office considers any general-purpose AI models that are

developed by the same entity through conducting another large pre-training run to constitute ‘distinct models’.

Recital 97 AI Act specifies that general-purpose AI models “*may be modified or fine-tuned into new models*”. The AI Office considers “fine-tuning” to be one way of “modifying” a general-purpose AI model. Currently, the AI Office considers modifications by the same entity to only lead to distinct models if those modifications use more than a third of the compute required for the model to having been classified as a general-purpose AI model (as specified in Section 3.1.1, i.e. currently roughly $3 \cdot 10^{21}$ FLOP) regarding the obligations for all providers of general-purpose AI models; and if those modifications lead to a significant change in systemic risk, which is presumed if those modifications use more than a third of the compute required for the model to having been classified as a general-purpose AI model with systemic risk (as specified in Article 51(2) AI Act, i.e. currently roughly $3 \cdot 10^{24}$ FLOP) regarding the obligations for providers of general-purpose AI models with systemic risk. By contrast, currently, the AI Office considers all other instances of a general-purpose AI model that are developed by the same entity and based on the same large pre-training run to constitute the same general-purpose AI model, and simply to be different ‘model versions’.

The differentiation between distinct models and model versions has implications on the actions providers of general-purpose AI models need to take to comply with their obligations:

- The documentation required under Article 53(1), points (a) and (b), AI Act must be separately drawn up for each distinct model placed on the market and updated when a different version thereof is made available.
- By contrast, the copyright policy required under Article 53(1), point (c), AI Act may be developed once by the provider and then applied to all its distinct models and versions thereof.
- The summary of the content used for training required under Article 53(1), point (d), AI Act must be separately drawn up and made publicly available for each distinct model placed on the market and updated as appropriate when a different version thereof is made available. The template that is currently under preparation by the AI Office will specify circumstances under which the template must be updated given a different model version.
- The systemic risk assessment and technical mitigations required under Article 55(1) AI Act must be carried out separately for each distinct model at appropriate milestones throughout its entire lifecycle. The development of a different version of a given distinct model may warrant a new systemic risk assessment and possibly technical mitigations. The general-purpose AI code of practice that is currently being drawn up will specify circumstances under which the development of a different version of a given distinct model requires a new systemic risk assessment and possibly technical mitigations. By contrast, the governance mitigations required under Article 55(1) AI Act may be developed once by the provider and then applied to all its distinct models and versions thereof.

The related question regarding the potential legal obligations faced by a downstream entity (distinct from the original provider) who modifies a given general-purpose AI model is addressed in Section 3.2.2.

3.2. Provider of a general-purpose AI model, including downstream modifiers

The definition of “provider” is key to understanding whether an entity must comply with the AI Act’s rules for general-purpose AI models.

Article 3(3) AI Act defines a “provider” of a general-purpose AI model as “*a natural or legal person, public authority, agency or other body that develops (...) a general-purpose AI model or that has (...) a general-purpose AI model developed and places it on the market (...)*”. According to Article 2(1), point (a), AI Act, providers of general-purpose AI models are in scope of the AI Act if they place a model on the Union market “*irrespective of whether they are established or located within the Union or in a third country*”. In the latter case, to facilitate compliance, providers have an obligation to “*appoint an authorised representative which is established in the Union*” prior to placing a model on the Union market, pursuant to Article 54 AI Act.

The definition of “provider” is not only relevant to the general-purpose AI part of the AI Act, but also to other parts of the AI Act relating to AI systems. To avoid a situation where important information is spread across several documents by the Commission, this section of the guidelines will focus specifically on clarifying who is considered to be subject to the obligations for providers of general-purpose AI models, by providing specific examples (Section 3.2.1), as well as by clarifying when a downstream entity may become the provider of a general-purpose AI model (Section 3.2.2).

3.2.1 Examples of providers of general-purpose AI models

The below list clarifies which entity is the provider of a general-purpose AI model in various hypothetical scenarios:

- If Entity A develops a general-purpose AI model and places it on the market, then Entity A is the provider.
- If Entity A has a general-purpose AI model developed for it by Entity B and Entity A places that model on the market, then Entity A is the provider.
- If Entity A develops a general-purpose AI model and uploads it to an online repository hosted by Entity C, then Entity A is the provider.
- If Entity A develops a general-purpose AI model and, on or outside the Union market, makes it available to Entity DM, who modifies the model in accordance with Section 3.2.2. and places the modified model on the market, then:
 - Entity A is the provider of the original model and must comply with the obligations for providers of general-purpose AI models, and
 - Entity DM is the provider of the modified model (“downstream modifier”, see Section 3.2.2) and must comply with the obligations for providers of general-purpose AI models, unless the modified model is not a general-purpose AI model.
- If Entity A develops a general-purpose AI model and, on or outside the Union market, makes it available to Entity DP who integrates the model into an AI system and makes the system available on the market or puts it into service, then:
 - Entity A is the provider of the model and must comply with the obligations for providers of general-purpose AI models, and
 - Entity DP is the provider of the system (“downstream provider”, see Article 3(66) AI Act) and must comply with the requirements for AI systems if applicable.

- If a collaborative or consortium has a general-purpose AI model developed for it by different individuals and/or organisations and places the model on the market, then usually the coordinator of the collaborative or the consortium is the provider. Alternatively, the collaborative or the consortium might be the provider. This must be assessed on a case-by-case basis.

3.2.2 Downstream modifiers as providers of general-purpose AI models

This section of the guidelines is expected to clarify responsibilities along the AI value chain, by specifying the conditions under which an entity who modifies a general-purpose AI model ('downstream modifier') must comply with the obligations for all providers of general-purpose AI models, and the obligations for providers of general-purpose AI models with systemic risk respectively, in the AI Act.

Recital 97 AI Act specifies that general-purpose AI models “*may be modified or fine-tuned into new models*”, and recital 109 AI Act further specifies that “*in the case of a modification or fine-tuning of a model, the obligations for providers of general-purpose AI models should be limited to that modification or fine-tuning, for example by complementing the already existing technical documentation with information on the modifications, including new training data sources, as a means to comply with the value chain obligations provided in this Regulation*”. The AI Office considers “fine-tuning” to be one way of “modifying” a general-purpose AI model.

At the same time, in line with the horizontal Blue Guide for New Legislative Frameworks acts ([Blue Guide](#)) such as the AI Act which states that “*a product which has been subject to important changes or overhauls aiming to modify its original performance, purpose or type may be considered as a new product*”, the AI Office deems that not every modification of a general-purpose AI model should lead to the downstream modifier being considered as a provider of a general-purpose AI model who is subject to the obligations laid down in Articles 52 to 55 AI Act. Instead, the AI Office deems that only those modifications that have a significant bearing on the rationales behind the obligations for providers of general-purpose AI models in the AI Act should lead to the downstream modifier being considered as a provider of a general-purpose AI model for the purposes of the respective obligations. For instance, when it comes to general-purpose AI models with systemic risk, only modifications that lead to a significant change in systemic risk should lead to downstream modifiers being considered as providers of general-purpose AI models with systemic risk.

In particular, the AI Office’s preliminary approach is to set certain thresholds in terms of computational resources used for the modification, that, if met, mean the downstream modifier should be presumed to be the provider of the modified general-purpose AI model or general-purpose AI model with systemic risk, and subject to the relevant obligations in the AI Act.

In this context, a distinction should be drawn between obligations for *all* providers of general-purpose AI models and obligations *only* for providers of general-purpose AI models with systemic risk. The former are laid down in Article 53 (excluding the obligation to document the information from Annex XI Section 2) and Article 54 AI Act. The latter are laid down in Article 52 and Article 55 AI Act.

The conditions and thresholds provided below may be updated by the AI Office in the future to reflect evolving technological developments regarding modifications and their impacts on model generality, capabilities, and systemic risk.

Downstream modifiers as providers of general-purpose AI models

A downstream modifier of a given general-purpose AI model who places the modified model on the market should be presumed to be the provider of a general-purpose AI model which has to comply with the obligations for all providers of general-purpose AI models with regards to the modification if the amount of computational resources used to modify the model is greater than a third of the training compute threshold for the original model as outlined in Section 3.1.1 (i.e. currently roughly $3 \cdot 10^{21}$ FLOP).

The indicative threshold of a third of the training compute threshold from Section 3.1.1 was chosen because when a model is modified using more than this amount of compute, the modified model can be expected to display a significant change in properties and behaviour, including a significant change in generality and capabilities as compared to the original model. These differences may be relevant for the AI Office, national competent authorities, and/or downstream providers, and justify the downstream modifier being subject to the obligations in Article 53(1), points (a) and (b), AI Act. A modification that meets this threshold can also be expected to have used a significant amount of data, which may be relevant for the copyright policy and the summary of the content used for training, justifying the modifier being subject to the obligations laid down in Article 53(1), points (c) and (d), AI Act. While as of now, few modifications may meet this threshold, more downstream modifiers may be in scope over time as the compute used to modify models increases.

In accordance with recital 109 AI Act, if a downstream modifier of a given general-purpose AI model becomes the provider of a general-purpose AI model based on the above criterion, then its obligations are limited to the modification conducted. That is, the documentation required by Article 53(1), points (a) and (b), AI Act need only concern the modification, and the copyright policy required by Article 53(1), point (c), AI Act as well as the summary of the content used for training required by Article 53(1), point (d), AI Act need only take into account the data used as part of the modification. The provider also has to comply with Article 54(1) AI Act.

Downstream modifiers as providers of general-purpose AI models with systemic risk

A downstream modifier of a given general-purpose AI model who places the modified model on the market should be presumed to be the provider of a general-purpose AI model with systemic risk which has to comply with the obligations for providers of general-purpose AI models with systemic risk if either of the following two conditions hold:

1. The original model is a general-purpose AI model with systemic risk, and the amount of computational resources used to modify the model is greater than a third of the training compute threshold for the original model as specified in Article 51(2) AI Act (i.e. currently roughly $3 \cdot 10^{24}$ FLOP);
2. The original model is not a general-purpose AI model with systemic risk, the downstream modifier knows or can reasonably be expected to know the cumulative amount of computational resources used to train this original model, and the sum of this amount and the amount of computational resources used to modify the model is greater

than the training compute threshold for the original model as specified in Article 51(2) AI Act (i.e. currently 10^{25} FLOP).

The threshold in point 1 has been chosen because, when a model is modified using a third of the training compute threshold for the original model to be considered a general-purpose AI model with systemic risk as specified in Article 51(2) AI Act (i.e. currently 10^{25} FLOP), the modified model can be expected to present significantly changed systemic risk compared to the original model. In particular, the AI Office deems that the original provider cannot be reasonably expected to have taken into account the change in systemic risk posed by such a modification in its systemic risk assessment and mitigation. As of today, the AI Office assumes that no downstream modification significantly changes systemic risk. The AI Office further assumes that, as of today, few or no modifications meet the specified threshold. The threshold in point 1 is thus forward-looking and in line with the risk-based approach of the AI Act.

The threshold in point 2 has been included to ensure that entities that modify general-purpose AI models that are not general-purpose AI models with systemic risk and which do not meet the threshold of the first condition, yet which lead to a modified model with cumulative training compute that is greater than the training compute threshold for the original model to be considered a general-purpose AI model with systemic risk specified in Article 51(2) AI Act (i.e. currently 10^{25} FLOP) are in scope of the obligations for general-purpose AI models with systemic risk. This is because in this case the original model would not have undergone any systemic risk assessment or mitigation by virtue of the model not being a general-purpose AI model with systemic risk, and therefore any systemic risks presented by the modified model will not have been assessed and mitigated by the provider of the original model.

If a downstream modifier of a given general-purpose AI model becomes the provider of a general-purpose AI model with systemic risk based on the above criterion, then its obligations are not limited to the modification conducted. That is, the systemic risk assessment and mitigation required by Article 55(1) AI Act should be conducted anew for the modified model, taking account any available information about the original model. The provider also has to notify the Commission in accordance with Article 52(1) AI Act.

3.3. Placing on the market of a general-purpose AI model and the open-source exemptions

The definition of “placing on the market” is key to understanding whether an entity must comply with the AI Act’s obligations for providers of general-purpose AI models.

Article 3(9) AI Act defines a ‘placing on the market’ of a general-purpose AI model as “*the first making available of (...) a general-purpose AI model on the Union market*”, while Article 3(10) AI Act defines a “making available on the market” of a general-purpose AI model as “*the supply of (...) a general-purpose AI model for distribution or use on the Union market in the course of a commercial activity, whether in return for payment or free of charge*”.

The guidelines are expected to build on the guidance already provided by recital 97 AI Act: “*(...) General-purpose AI models may be placed on the market in various ways, including through libraries, application programming interfaces (APIs), as direct download, or as physical copy. (...) It should be understood that the obligations for the providers of general-purpose AI models should apply once the general-purpose AI models are placed on the market. When the provider of a general-purpose AI model integrates an own model into its own AI system that is made available on the market or put into service, that model should be considered*

to be placed on the market and, therefore, the obligations in this Regulation for models should continue to apply in addition to those for AI systems. The obligations laid down for models should in any case not apply when an own model is used for purely internal processes that are not essential for providing a product or a service to third parties and the rights of natural persons are not affected. Considering their potential significantly negative effects, the general-purpose AI models with systemic risk should always be subject to the relevant obligations under this Regulation. (...)”

The definition of “placing on the market” is not only relevant to the general-purpose AI part of the AI Act, but also to other parts of the AI Act relating to AI systems. To avoid a situation in which important information is spread across several documents by the Commission, these guidelines will focus specifically on clarifying what is considered to be a placing on the market of a general-purpose AI model, by providing specific examples (Section 3.3.1), as well as by clarifying when providers of general-purpose AI models released as open-source benefit from the exemptions from certain obligations (Section 3.3.2).

3.3.1 Examples of placing on the market of general-purpose AI models

The below list provides various examples of when a general-purpose AI model should be considered to have been placed on the market:

- A general-purpose AI model is made available via a software library or package
- A general-purpose AI model is made available via an application programming interface (API)
- A general-purpose AI model is uploaded to a public repository for direct download (for details, see Section 3.3.2)
- A general-purpose AI model is made available as a physical copy
- A general-purpose AI model is made available via a cloud computing service
- A general-purpose AI model is copied onto a customer’s own infrastructure
- A general-purpose AI model is integrated into a chatbot made available via a web interface
- A general-purpose AI model is integrated into a mobile application made available through app stores
- A general-purpose AI model is integrated into the provider’s own products or services that are made available on the market

These examples should be interpreted in accordance with the horizontal guidance for the definition of ‘placing on the market’ as provided for in the Blue Guide and the relevant provisions in the AI Act (Article 3(9) and (10), and recital 97 AI Act).

3.3.2 Exemptions from certain obligations for certain open-source releases

This section of the guidelines is expected to clarify when the exemptions from certain obligations for providers of general-purpose AI models released as open-source apply and, in particular, the conditions under which these exemptions apply.

The AI Act applies, in principle, to providers of general-purpose AI models placed on the market as open-source. However, according to recital 102 AI Act, “[s]oftware and data, including models, released under a free and open-source licence that allows them to be openly

shared and where users can freely access, use, modify and redistribute them or modified versions thereof, can contribute to research and innovation in the market and can provide significant growth opportunities for the Union economy.” Therefore, the AI Act foresees exemptions from several obligations for providers of such models meeting specific conditions.

Firstly, Article 53(2) AI Act contains an exemption from the obligations to document information for the purposes of providing that information, upon request, to the AI Office and national competent authorities, and making it available to downstream providers who intend to integrate the model into their systems: *“The obligations set out in paragraph 1, points (a) and (b), shall not apply to providers of AI models that are released under a free and open-source licence that allows for the access, usage, modification, and distribution of the model, and whose parameters, including the weights, the information on the model architecture, and the information on model usage, are made publicly available. This exception shall not apply to general-purpose AI models with systemic risks.”* Secondly, Article 54(6) AI Act contains an exemption from the obligation to appoint an authorised representative for providers of general-purpose AI models established in third countries under the same conditions: *“The obligation set out in this Article shall not apply to providers of general-purpose AI models that are released under a free and open-source licence that allows for the access, usage, modification, and distribution of the model, and whose parameters, including the weights, the information on the model architecture, and the information on model usage, are made publicly available, unless the general-purpose AI models present systemic risks.”*

To benefit from the exemptions, the provider’s model must meet the following conditions:

- **The general-purpose AI model is released under a free and open-source licence that allows for the access, usage, modification, and distribution of the model:**
 - “Access” means that the licence foresees that anybody interested can freely obtain the model without any payment requirements or other restrictions.
 - “Usage” means that the licence guarantees that the original provider will not use their intellectual property rights to restrict, or charge for, the use of the model, subject to limited conditions. Such limited conditions may only serve to ensure that *“the original provider of the model is credited the identical, [and] comparable terms of distribution are respected”* (recital 102 AI Act).
 - “Modification” means that the licence allows anybody to freely make alterations to the model without any payment requirements or other restrictions.
 - “Distribution” means that the licence allows those who access, use and modify the general-purpose AI model to freely distribute it onwards, subject to limited conditions. Again, such limited conditions may only serve to ensure that *“the original provider of the model is credited, [and that] the identical or comparable terms of distribution are respected”* (recital 102 AI Act).
 - “Free and open-source” includes not only the concepts of free access, use, modification and distribution as explained above, but also indicates that the model, along with its associated services, should be provided at no monetary compensation (recital 103 AI Act). For the exemption to apply there should not be any monetary compensation demanded in exchange for access to the general-purpose AI model or associated services. Monetisation includes provision of the model against a price or any other monetisation. For transactions other than transactions between microenterprises, monetisation includes provision of technical support and other services, or the use of personal data collected from the use of the model or the accompanying services, except when the personal

data is used exclusively to improve the security, compatibility or interoperability of the software. “*The fact of making [general-purpose AI models] available through open repositories [does] not, in itself, constitute a monetisation.*” (recital 103 AI Act).

- **The parameters, including the weights, the information on the model architecture, and the information on model usage, are made publicly available:** this requires that the files containing model parameters, including the weights, as well as supporting documentation are provided in a format, and with a degree of clarity and specificity, that enable usage, modification and distribution in accordance with the free and open-source licence.
- **The general-purpose AI model is not a general-purpose AI model with systemic risk.**

3.4. Estimating the computational resources used to train or modify a model

The AI Act includes a threshold involving computational resources (compute) used to train a general-purpose AI model which can lead the model to be classified as a general-purpose AI model with systemic risk. Sections 3.1.1 and 3.2.2 introduce thresholds which also involve compute. To know whether a model meets any of these thresholds, potential providers must estimate the amount of compute used. This section of the guidelines will aim to provide guidance on how they may do so. Section 3.4.1. describes two widely used approaches that potential providers may use to estimate the compute used for training or modifying an AI model. Section 3.4.2 specifies how potential providers may further estimate the *cumulative* amount of compute used to train their general-purpose AI model, in light of Articles 51(2), 52(1) and recital 111 AI Act, including what should be included in the calculation, and how and when it should be estimated.

3.4.1. Estimating the amount of compute used for training or modification

The amount of compute used to train or modify a model can be estimated in two ways: by tracking Graphics Processing Unit (GPU) usage (hardware-based approach), or by counting the expected number of FLOP based on the model’s architecture (architecture-based approach). Potential providers may choose which of the two approaches to use to determine whether their model meets the thresholds set out in Section 3.1.1 or 3.2.2 respectively. In either case, all operations should be counted equally, regardless of floating-point precision. Moreover, approximations that modify the final result by less than 5% are presumed to be valid.

Hardware-based approach

To estimate their training compute based on hardware, the potential provider of a general-purpose AI model should first identify:

- The number N of GPUs used for the training or modification;
- The total training duration L (measured in seconds);
- The peak theoretical performance H of the GPUs used (measured in FLOP/second, calculated via weighted average if different types of GPUs are used and/or the same types of GPUs are used with different number formats);
- The percentage of GPU utilisation U achieved.

The amount C in FLOP of training compute can then be calculated according to the following formula:

$$C = N \cdot L \cdot H \cdot U.$$

Section A.1. of the Annex includes example calculations using this approach.

Architecture-based approach

This approach consists in estimating the computational resources used for training or modification based on the model's architecture.

AI models based on neural networks (which cover essentially all general-purpose AI models or potentially general-purpose AI models today) are trained through a succession of forward and backward passes. The amount C in FLOP of training compute is the product of the number of full passes (one full pass is the combination of a forward and backward pass) made during training with the total number of operations performed in a full pass, i.e.

$$C = \text{Number of forward passes made during training} \cdot \text{number of operations/full pass.}$$

For large⁴ models based on the transformer architecture, providers may use the following approximation for estimating C :

$$C \approx 6 \cdot P \cdot D,$$

where P refers to the total number of model parameters active per forward pass and D refers to the total number of training examples⁵ used for training (see for example <https://arxiv.org/abs/2001.08361>).

Section A.1. of the Annex includes example calculations using this approach.

3.4.2. Estimating the cumulative amount of computational resources used for training

What should be counted?

According to Article 51(1), point (a), AI Act, a general-purpose AI model is classified as a general-purpose AI model with systemic risk if it has “*high impact capabilities evaluated on the basis of appropriate technical tools and methodologies, including indicators and benchmarks*”. In addition, Article 51(2) AI Act provides that “*a general-purpose AI model shall be presumed to have high impact capabilities pursuant to paragraph 1, point (a), when the cumulative amount of computation used for its training measured in floating point operations is greater than 10^{25} .*”

Regarding how the “cumulative amount of computation” used to train a model (cumulative training compute) should be understood, recital 111 AI Act specifies that “*the cumulative amount of computation used for training includes the computation used across the activities and methods that are intended to enhance the capabilities of the model prior to deployment, such as pre-training, synthetic data generation and fine-tuning*”.

The AI Office understands such activities and methods as being restricted to activities and methods carried out as part of the training of the model (for example pre-training, fine-tuning, Reinforcement Learning from Human Feedback), or directly feeding into the training of the model (for example synthetic data generation), and therefore not to include activities carried out prior to the large pre-training run (for example research and development activities,

⁴ A model based on a transformer architecture should be considered large if its number of parameters is significantly larger than one twelfth of the number of tokens in the input context. This ensures that the contribution to the training compute coming from the context-dependent compute is negligible in comparison to the contribution coming from the non-embedding compute and which equals $6PD$ (see for example <https://arxiv.org/abs/2001.08361>).

⁵ Where the training data is text, a training example corresponds to a training token..

prototyping and testing, smaller test runs, and hyperparameter tuning), or which improve the model's capabilities at inference time (for example scaffolding).

How should it be counted?

To estimate the cumulate training compute of their model, providers may use either the hardware-based approach or the architecture-based approach specified in Section 3.4.1, ensuring in either case that all activities and methods carried out as part of training or directly feeding into training, and intended to enhance model capabilities, have been accounted for in line with the above guidance.

To account for the compute used to generate synthetic data, the AI Office recognises that in cases where the provider has generated this data from model(s) that it has not developed itself, or where the provider has obtained the synthetic data set(s) from a third party, the provider may not be able to directly calculate the amount of compute used to generate the synthetic data. This is further complicated by the fact that in the case of synthetic data generation, multiple outputs may be generated and filtered to produce a single high-quality input (for example via rejection sampling), and providers may not know the rejection rate. For this reason, the AI Office's preliminary approach is to allow providers to use reasonable estimates when precise information is impractical to obtain. In this case, providers should document their method of estimation. This approach is likely to change in the light of evolving technological developments.

Where a model has been created by combining a number of smaller models, for example via the Mixture-of-Experts technique, or through integrating pre-existing model parameters, for example through parameter initialisation, the training compute used to train the original models should be included in the estimation of the cumulative training compute of the final model.

When should it be counted?

This Section of the guidelines will clarify when a provider should notify the Commission that their general-purpose AI model has or will meet the compute threshold from Article 51(3) AI Act, in accordance with Article 52(1) AI Act.

According to Article 52(1) AI Act, the provider of a “general-purpose AI model that meets the condition referred to in Article 51(1), point (a),” AI Act must “notify the Commission without delay and in any event within two weeks after that requirement is met or it becomes known that it will be met.” Moreover, “the notification shall include the information necessary to demonstrate that the relevant requirement has been met.” Regarding when a provider is expected to know that their general-purpose AI model will meet the threshold specified in Article 51(2) AI Act, recital 112 AI Act specifies that “training of general-purpose AI models takes considerable planning which includes the upfront allocation of compute resources and, therefore, providers of general-purpose AI models are able to know if their model would meet the threshold before the training is completed.”

The AI Office assumes that “planning” and “upfront allocation of compute resources” referred to in recital 112 AI Act are activities that take place prior to the start of the large pre-training run, and that they allow for estimation of the amount of pre-training compute that will be used to train the model. Whilst the cumulative amount of compute may be larger still (though within the current paradigm for training the most advanced general-purpose AI models, pre-training compute still represents the vast majority of the cumulative amount of training compute), providers may not know before they begin pre-training what this amount will be.

Accordingly, the AI Office’s preliminary approach is to expect providers to estimate the amount of *pre-training* compute that they will use ahead of commencing their large pre-training run (using either of the two methods specified in Section 3.4.1), and to notify the Commission “without delay and in any event within two weeks” if the estimated value meet the threshold specified in Article 51(2) AI Act, following Article 52(1) AI Act.

If the estimated value does not meet the threshold specified in Article 51(2) AI Act, providers are expected to closely monitor their actual and expected compute usage over the course of training, including after the large pre-training run is completed, so that they are able to know if and when the cumulative amount of compute has met or will meet the threshold, and notify the Commission accordingly following Article 52(1) AI Act.

When the provider of a general-purpose AI model knows that their model has met or will meet the training compute threshold, the AI Office interprets “*the information necessary to demonstrate that the relevant requirement has been met*” which the provider has to include with their notification as information which should at minimum contain:

- The amount of compute estimated by the provider which has triggered the requirement to notify, reported in FLOP and with two significant figures; and
- A description of the approach used to estimate this amount of compute, including approaches used for making approximations where precise information is not readily available.

3.5. Transitional rules, grandfathering, and retroactive compliance

According to Article 111(3) AI Act, “[p]roviders of general-purpose AI models that have been placed on the market before 2 August 2025 shall take the necessary steps in order to comply with the obligations laid down in this Regulation by 2 August 2027.”

The AI Office recognises that in the months following the entry into application of the obligations of providers of general-purpose AI models in the AI Act on 2 August 2025, some providers may face various challenging situations to ensure timely compliance with their obligations under the AI Act. Accordingly, the AI Office is dedicated to supporting providers in taking the necessary steps to comply with their obligations. In particular:

- For general-purpose AI models that have been placed on the market before 2 August 2025, providers must take the necessary steps to comply with their obligations by 2 August 2027. This does not require re-training or unlearning of models already trained before 2 August 2025, where implementation of the measures for copyright compliance is not possible for actions performed in the past, where some of the information for the training data is not available, or where its retrieval would cause the provider disproportionate burden. Such instances must be clearly justified and disclosed in the copyright policy and the summary of the content used for training.
- For general-purpose AI models with systemic risk that have been placed on the market before 2 August 2025, providers who foresee difficulties with complying with their obligations should reach out to the AI Office to discuss how they can be supported in complying with their obligations.
- For providers who, on 2 August 2025, have trained, are in the process of training, or are planning on training a general-purpose AI model with a view to placing the model on the market after 2 August 2025, and who foresee difficulties with complying with the obligations for all providers of general-purpose AI models, they should proactively

inform the AI Office how and when they will take the necessary steps to comply with their obligations.

- For providers who, on 2 August 2025, have trained, are in the process of training, or are planning on training a general-purpose AI model with systemic risk with a view to placing the model on the market after 2 August 2025, the AI Office expects providers to notify it of the model without delay and in any event within two weeks after 2 August 2025 as per Article 52(1) AI Act. If providers foresee difficulties with complying with the other obligations for providers of general-purpose AI models with systemic risk, they should flag those to the AI Office for consideration alongside their notification. In the special case where a provider has never placed on the market a general-purpose AI model with systemic risk before 2 August 2025, the AI Office will give particular consideration to their challenging situation with respect to setting any deadlines for taking the necessary steps to comply with their obligations to allow a timely placement on the market.

3.6. Effects of signature and adherence to a code of practice

Providers of general-purpose AI models that are signatories to a code of practice will be transparent in their compliance with the AI Act and therefore benefit from increased trust by the Commission and other stakeholders. While compliance with the AI Act can be demonstrated through various means, adherence to a code of practice approved by the AI Office and the Board is a straightforward and transparent means to demonstrate compliance with the AI Act (Articles 53(4) and 55(2) AI Act). For signatories to a code of practice, the Commission can be expected to focus its enforcement on monitoring adherence to the code of practice (Article 89(1) AI Act). The Commission may approve a code of practice via implementing act, thereby giving it a general validity within the Union (Article 56(6) AI Act).

Non-signatories will be expected to demonstrate how they comply with their obligations under the AI Act via other adequate, effective, and proportionate means by reporting the measures they have taken to the AI Office. Furthermore, non-signatories will be expected to explain how the measures they have taken ensure compliance with their obligations under the AI Act, for instance by carrying out a gap analysis. Finally, non-signatories may be subject to more requests for information and access to conduct model evaluations, since there may be less clarity regarding how they ensure compliance with their obligations under the AI Act.

The Commission may take into account commitments made in a code of practice as a mitigating factor when fixing the amount of fines, depending on the specific circumstances (Article 101(1) AI Act).

Providers of general-purpose AI models without systemic risk may voluntarily choose to adhere to commitments relevant for providers of general-purpose AI models with systemic risk (Article 56(7) AI Act) in relation to their general-purpose AI models without systemic risk. In this case, the AI Office will only monitor adherence to the commitments relevant for general-purpose AI models with systemic risk once the provider develops a general-purpose AI model with systemic risk (Article 89(1) AI Act).

Commitments made in a code of practice become relevant for assessing the provider's compliance with the AI Act only with the entry into application of the obligations of providers of general-purpose models in the AI Act on 2 August 2025.

If a code of practice cannot be finalised by 2 August 2025, the Commission may adopt common rules via implementing act, which would be applicable to all providers of general-purpose AI models and general-purpose AI models with systemic risk (Article 56(9) AI Act). A code of practice is a temporary tool until harmonised standards are developed (Article 40 AI Act). If, after a standardisation request, harmonised standards are not developed in time or in a satisfactory manner, the Commission may adopt common specifications via implementing act (Article 41 AI Act).

3.7. Supervision and enforcement of the general-purpose AI rules

The AI Office will supervise and enforce the obligations laid down in the AI Act for providers of general-purpose AI models (Article 88 AI Act) and the compliance of the AI systems based on general purpose AI models if the provider of the model and the system are the same (Article 75(1) AI Act). The following clarifications concern the supervision and enforcement of the obligations for providers of general-purpose AI models.

The AI Office will take a collaborative and proportionate approach to enforcement. The AI Office expects close informal cooperation with providers during the training of the general-purpose AI model to streamline compliance and ensure market placement without delays, in particular for providers of general-purpose AI models with systemic risk. Furthermore, the AI Office expects proactive reporting without requests by providers of general-purpose AI models with systemic risk, as part of either commitments under a code of practice or alternative means to demonstrate compliance.

Enforcement by the AI Office is underpinned by the powers given to it under the AI Act, namely the powers to request information (Article 91 AI Act), conduct evaluations of general-purpose AI models (Article 92 AI Act), request measures from providers, including implementing risk mitigations and recalling the model from the market (Article 93 AI Act), and to impose fines of up to 3% of global annual turnover or EUR 15 million, whichever is higher (Article 101 AI Act) starting on 2 August 2026 after a grace period of one year. More detailed acts will follow to further specify the implementation of these powers, in particular the implementing acts derived from Articles 92 and 101 AI Act.

The AI Office will ensure the confidentiality of the data obtained in carrying out its tasks and activities in order to protect, in particular, intellectual property rights, confidential business information or trade secrets of natural or legal persons, and public safety and security, in accordance with Article 78 AI Act.

Annex

A.1. Compute used for training of models with approximately one billion parameters

The preliminary proposal for the threshold for a model to qualify as a general-purpose AI model proposed in Section 3.1.1 has been informed by the estimated amount of compute used in practice to train widely used models with approximately one billion parameters. While we do not refer to the models and their providers by name, we share the information relevant to our preliminary estimate in the following. The estimate is based on four models:

- Model A (a language model with 1.1 billion parameters): $2 \cdot 10^{22}$ FLOP.
- Model B (a language model with 1 billion parameters): $3 \cdot 10^{21}$ FLOP.
- Model C (a language model with 3.8 billion parameters): $7.5 \cdot 10^{22}$ FLOP.
- Model D (an image diffusion model with 1 billion parameters): $1.2 \cdot 10^{23}$ FLOP.

The above values have been obtained via the following calculations, which also serve to illustrate the two estimation methods described in Section 3.3.1.

Model A:

Architecture-based approach: Model A is a language model based on a transformer decode architecture with 1.1 billion parameters, trained on 3 trillion tokens. Applying the approximate formula for the training compute C from Section 3.4.1 yields

$$C \approx 6PD = 6 \cdot 1.1 \cdot 10^9 \cdot 3 \cdot 10^{12} = 1.98 \cdot 10^{22} \text{ FLOP.}$$

Hardware-based approach: Model A's documentation also mentions that it took 3456 GPU hours to train it on 300B tokens using A100 GPUs, which implies 35,000 GPU hours for the full training run given the model was trained on 3 trillion tokens. This implies:

- a product NL of number of GPUs N and total training duration of L equal to $NL = 35,000 \cdot 3,600$ seconds;
- a peak performance per GPU of $H = 300 \cdot 10^{12}$ FLOP/second as per the NVIDIA A100 datasheet;⁶
- a GPU utilization of $U = 56\%$ as indicated in model's documentation.

Using the hardware-based approach in Section 3.4.1, the total compute is then given by

$$C = N \cdot L \cdot H \cdot U \approx 2.1 \cdot 10^{22} \text{ FLOP,}$$

confirming the approximate estimate obtained via the architecture-based approach.

Model B:

Hardware-based approach: Model B is a language model with a transformer decoder architecture with one billion parameters, trained with 4830 GPU hours on 40GB A100 GPUs. Assuming the same numbers for peak performance and GPU utilisation as for Model A, this implies a total compute of

⁶ <https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/a100/pdf/nvidia-a100-datasheet-us-nvidia-1758950-r4-web.pdf> assuming bfloat16 training with no sparsity.

$$C = N \cdot L \cdot H \cdot U \approx 3 \cdot 10^{21} \text{ FLOP.}$$

Model C:

Architecture-based approach: Model C is a language model with a transformer decode architecture and 3.8 billion parameters, trained on 3.3 trillion tokens. The approximation from Section 3.4.1 yields

$$C \approx 6 \cdot P \cdot D = 6 \cdot 3.8 \cdot 10^9 \cdot 3.3 \cdot 10^{12} \approx 7.5 \cdot 10^{22} \text{ FLOP.}$$

Model D:

Hardware-based approach: Model D is a well-known diffusion model for image generation. According to its model card it was trained for a total of 200,000 GPU hours on 40GB A100 GPUs. Assuming the same numbers for peak performance and GPU utilisation as previously, this implies a total compute of

$$C = N \cdot L \cdot H \cdot U \approx 1.2 \cdot 10^{23} \text{ FLOP.}$$