# DEVELOPMENT OF IA SYSTEMS: WHAT SHOULD BE CHECKED?

The CNIL offers a **list of points to check**, taken from its **recommendations for the development of RGPD-compliant AI systems.**

It is aimed at designers and developers of artificial intelligence systems (product managers, developers, data protection officers, legal teams, information systems security officers, etc.) in order to secure all the stages in the development of an AI system, **from data collection to integration**, including **model learning** and **annotation**. This checklist is designed to ensure that, from the outset, the principles of the RGPD are correctly implemented: purpose, minimisation, security, information, people's rights, transparency and governance.

Please note: the obligations laid down by the Regulation on artificial intelligence must, where applicable, also be taken into account in the development of these systems. They are not covered by this checklist.

| SHEETS | | MEASURES | |
|---|---|---|---|
| **1** Determine the applicable legal regime and your liability Define a framework for users | Identify whether the RGPD applies | Identify whether the training database contains **personal data** (including from *web scraping*). | ☐ |
| | | Analyse whether the RGPD applies to the model learned from training databases c o n t a i n i n g personal data or whether it can be presumed to be anonymous. • To do this, establish whether it is necessary to carry out re-identification attacks on the AI model, the profundity of the model and the anonymity of the model. of these attacks, and the likelihood of personal data being extracted, detailing as much as possible by type of data. | ☐ |
| | | If you consider that the system incorporating a non-anonymous AI model may make it possible to escape the scope o f the RGPD, check that the measures put in place are sufficiently effective and robust to make the likelihood of re-identification of individuals insignificant. • This assessment necessarily involves carrying out re-identification attacks on the AI system. | ☐ |
| | | Implement a process for regular reassessment of the anonymity of the model or system. | ☐ |
| | Defining the responsibilities of those involved | Determine your **responsibility** and that of other parties involved in the processing of personal data- (controller, joint controller or processor). | ☐ |
| | | Where appropriate, ensure that you have a contract to **govern joint responsibility.** | ☐ |
| | | If applicable, ensure that you have a contract to **document the instructions** given to your processors. | ☐ |
| **2** Define the purposes and choose the legal basis | Define the purposes | Clarify **the purpose(s) of the project** at the design stage. • In the case of general-purpose AI, refer to the type of system being developed (e.g. the development o f a large-scale language model, a computer vision system) and the technically feasible functions and capabilities. | ☐ |
| | Identify the legal basis | Identify **the legal basis(s)** for each processing operation (consent, legitimate interest, etc.). | ☐ |
| | | Where appropriate, document **the procedures for obtaining consent** and keep proof of this (article 6.1.a of the RGPD). | ☐ |
| | | Where applicable, ensure that you have a **valid contract** and that the processing is necessary to fulfil the purpose of the contract (article 6.1.b of the RGPD). | ☐ |
| **3** If necessary, assess the validity of the legal basis of the legitimate interest (Article 6.1.f of the RGPD). | Check the existence of a legitimate interest | Clearly define **the interest pursued.** | ☐ |
| | | Check that the interest does not **conflict with other regulatory obligations** (Digital Services Regulation, Artificial Intelligence Regulation, etc.). | ☐ |
| | Assess the need for processing | Check that the processing of personal data is **necessary to achieve the defined objective.** | ☐ |
| | | Check that **less intrusive methods** (e.g. anonymisation, synthetic data) cannot achieve the same results. | ☐ |
| | | Check that the algorithmic techniques used to process the data (e.g. convolutional deep neural networks, SVM support vector machines, etc.) **consume** the **least amount of personal data** possible for the objective pursued. Where appropriate, document the need to use machine learning, in particular deep learning. | ☐ |
| | | Check whether design choices can be taken into account with a view to protecting data right from the design stage (federated learning, secure multi-party computing, homomorphic encryption, etc.). | ☐ |
| | Weighing up the interests at stake | Ensure and document that data subjects can **reasonably expect** this processing to be carried out. | ☐ |
| | | Where appropriate, implement and **document suitable and sufficient safeguards to limit the impact of the processing on the data subjects** (e.g. provide for the rapid anonymisation of the data collected or, failing that, the pseudonymisation of the data collected, adopt measures to limit the risks of storage, extraction, regurgitation in the context of generative AI or attacks on AI models or systems, provide for a discretionary and prior right to object, etc.). | ☐ |

**CNIL.**

| | | | | |
|---|---|---|---|---|
| | | | Where necessary, implement **appropriate safeguards against** *web scraping of* data (e.g. limit data collection to freely accessible data, draw up a list of sites from which collection would be excluded by default because they contain particularly intrusive data). | ☐ |

| SHEETS | | | MEASURES | |
|---|---|---|---|---|
| 4 | When re-using data, carry out additional tests and checks | If you re-use your own data | If the training purpose of your model was not foreseen at the time of data collection, check that it is compatible with the initial purpose using the **compatibility test** (unless you are authorised by the data subjects because they have given their consent or by a text, or if you are re-using the data for statistical or scientific research purposes):<br>• Is there a link between the initial purpose and the new AI purpose?<br>• Does the context of the initial collection reasonably allow this re-use?<br>• What is the type and nature of the data (identifiers, sensitive, etc.)?<br>• What are the possible consequences for individuals?<br>• What technical and organisational safeguards are in place (pseudonymisation, etc.)? | ☐ |
| | | If you are re-using publicly accessible data or data acquired from a third party (e.g. data broker) | Check that you are not re-using a **database whose creation was manifestly illegal:**<br>• Is the source of the data clearly identified and documented?<br>• Is the source of the data clearly identified and documented? Is the database not clearly the result of a crime or misdemeanour (leakage, theft, etc.) or has it been the subject of a conviction or public sanction by a competent authority that has led to it being suppressed or banned from use?<br>• Are the conditions under which the data is collected sufficiently documented?<br>• Does the database not contain sensitive or infringement data, or are enhanced checks made to ensure that the processing is lawful, if so? | ☐ |
| 5 | Limit the data processed to that which i s  relevant and necessary (data minimisation) | Selection of strictly necessary data | Identify the data that is essential to achieving your purposes and favour less intrusive formats (e.g. age range rather than full date of birth). | ☐ |
| | | | **As regards the volume of data**, justify the number of people concerned, the historical depth and the seriousness. | ☐ |
| | | | Justify the need to process **highly personal data.** | ☐ |
| | | | **With regard to the type of data**, assess the use of real data and synthetic, pseudonymised or anonymised data. | ☐ |
| | | | Identify **data sources.** | ☐ |
| | | Implement specific measures in the event of *web scraping*. | Define **precise** collection **criteria** beforehand. | ☐ |
| | | | Exclude the collection of **certain categories of data** when they are not necessary, using filters where possible or, failing that, by excluding certain types of site that structurally contain these categories of data. | ☐ |
| | | | Exclude from collection **sites that clearly oppose harvesting** of their content, for example by using robots.txt files or CAPTCHA. | ☐ |
| | | Specific precautions for sensitive data | Justify the **need** to process sensitive data. | ☐ |
| | | | Identify **the exception to the principle of prohibiting the processing of sensitive data** (Article 9.2 of the RGPD). | ☐ |
| | | | Provide for **enhanced security** measures (pseudonymisation, etc.). | ☐ |
| | | | Immediately and, if possible, automatically delete sensitive data collected incidentally or- during data harvesting (*web scraping*). | ☐ |
| | | Organisation of data collection and preparation | Carry out data **cleansing** (inconsistencies, duplicates, etc.). | ☐ |
| | | | Identify the **data that is really relevant to** the task and delete data that is not relevant for learning. | ☐ |
| | | | Apply **protection techniques at the design stage** (e.g. generalisation, randomisation, pseudonymisation, anonymisation, etc.). | ☐ |
| | | Justification and validation of design choices | Carry out a **pilot experiment** with fictitious, synthetic or anonymised data. | ☐ |
| | | | Ask a **referee or an ethics committee** about the ethical issues and the protection of people's rights rights and freedoms. | ☐ |
| | | Ongoing review | Implement a process for regularly reviewing the relevance of the data collected. | ☐ |
| | | | Implement mechanisms for deleting unnecessary or obsolete data. | ☐ |
| 6 | Define and control data retention periods | Define a clear retention policy at the design stage | Define a specific retention period for each phase of the IA project lifecycle (development, maintenance, enhancement, etc.). | ☐ |
| | | Retention during the development phase | Check that data is only accessible to authorised personnel during development. | ☐ |
| | | | Establish a process for archiving or deleting data at the end of the development phase, unless it is necessary to retain the data for maintenance or product improvement. | ☐ |
| | | | Document the need to retain data beyond the development phase, in particular for product maintenance or improvement. | ☐ |
| | | Retention for the | Check that data is stored on a partitioned, secure medium. | ☐ |

CNIL.

| | | | | |
|---|---|---|---|---|
| **9** | | purposes of product maintenance or improvement | Check that access is strictly limited to the people in charge of maintaining or improving the the product. | ☐ |
| | | | Implement an automatic deletion plan once improvements have been made. | ☐ |

| **SHEETS** | | | **MEASURES** | |
|---|---|---|---|---|
| **7** | Ensure transparency of processing operations | | **Inform** data subjects in a clear and easily accessible manner of all the information provided for in Articles 13 and 14 of the RGPD. | ☐ |
| | | | Where data is not collected directly from individuals, document, where applicable, the fact that individual information **would require disproportionate efforts** (article 14.5.b of the RGPD) and make the information publicly available (website, etc.). | ☐ |
| | | | If *web scraping* concerns a limited number of sites, **provide precise information** on the sources used. If the sources are very numerous, **provide the categories of source sites**, at least those presenting the greatest risk to individuals. | ☐ |
| | | | Where the model itself is subject to the GDPR, **inform people about the data stored** and the risks to **which they are exposed.** provide all the information required by Articles 13 and 14 of the RGPD. | ☐ |
| **8** | Respecting people's rights | Establishing procedures for managing rights | Inform individuals of the **risk of data regurgitation** in the case of generative AI, the **measures taken to limit these risks** and the **existing recourse mechanisms** (e.g. the possibility of reporting an occurrence of regurgitation to the or-ganism) | ☐ |
| | | | Establish a **procedure for notifying recipients, in particular users, of any request** for rectification, erasure or restriction of data, unless such notification proves impossible or would require disproportionate effort. | ☐ |
| | | Manage the identification of individuals | In the case of generative AI, **establish an internal procedure consisting of querying the model** (for example from a list of selected queries) to check the data it may have stored on the person thanks to the information provided. | ☐ |
| | | | If it is not possible to identify a person within the model but they are identified in a training d a t a b a s e , inform them of the risk of memorising the model. | ☐ |
| | | | If it is **not possible to identify a person** in the training database or in the model, inform them of this. inform them of this. | ☐ |
| | | | Inform the person of **the additional information that** can be provided to help identify them (e.g. pseudonym, sample of their data). | ☐ |
| | | | Implement a process for **deleting this additional data** once the request has been processed. | ☐ |
| | | Choose a technical solution to ensure that model rights are respected | In principle, provide for **a** model **retraining process**.<br>• Re-training can be periodic to limit costs and satisfy several requests to exercise rights at the same time. | ☐ |
| | | | Provide users with an **updated version of the model**, possibly by contractually requiring them to use only a regularly updated version. | ☐ |
| | | | Where applicable, document the fact that **re-training of the model proves disproportionate** (temporarily or permanently). | ☐ |
| | | | If the re-training is disproportionate, **implement filters or other robust measures** on the output of the AI system. | ☐ |
| | | | Where appropriate, **prefer general rules** preventing the generation of personal data to a simple "blacklist" of people who have exercised their rights. | ☐ |
| **9** | Ensuring data annotation i s compliant | | Check that annotations contain only the information necessary to achieve the purpose and that they that they are objective. | ☐ |
| | | | Implement **regular reviews** to ensure the continued relevance of labels. | ☐ |
| | | | Implement a **continuous verification procedure** to control the quality of annotation: define an annotation protocol, applying the principles of accuracy and minimisation, and involve a referee or an ethical committee (good practice). | ☐ |
| | | | Inform data subjects of the data annotation phase. | ☐ |
| | | | Ensure that **internal procedures for managing rights and the methods for exercising rights** include annotation (right of access, rectification, erasure, restriction, portability, opposition). | ☐ |
| | | | Check, where applicable, that sensitive data is processed in accordance with an exception to the principle of the prohibition of data (article 9.2 of the RGPD). | ☐ |

CNIL.

| | | | | |
|---|---|---|---|---|
| | | | Implement **specific measures** in view of the increased risk to individuals: annotate according t o objective and factual criteria, limit annotation to the context of the data, strengthen the annotation verification stage, increase the security of annotated data (eg. increase the security of the annotated data (e.g. by carrying out the annotation processing in-house, by processing the data locally and by guaranteeing its security through encryption, logging and stronger access restrictions) and consider the risk of regurgitation and inference of sensitive data on the models trained from it. | ☐ |
| | | | Train the people in charge of annotation in data protection principles. | ☐ |
| **SHEETS** | | | **MEASURES** | |
| 10 | Ensure the security of data | | Check that the security measures relating to training data are adequate and appropriate (see the guide to personal data security). | ☐ |
| | | | Check that the security measures relating to the development of the system are sufficient and appropriate. In particular, use verified development tools, libraries and, where appropriate, pre-trained models. | ☐ |
| | | | Check that the measures designed to control the operation of the system are sufficient and appropriate. For example: favour verified import and backup formats such as *safetensors*, control the output of the AI system using filters, implement *watermarking* techniques. | ☐ |
| | | | In general terms, managing access rights to data, tracing access and analysing traces. | ☐ |
| | | | Implement and monitor an action plan to ensure that safety requirements are met. | ☐ |
| 11 | Analyse the risks and carry out a data protection impact assessment (DPIA) | | Carry out a DPIA if the model training process presents **high risks** based on the criteria identified by the European Data Protection Committee (innovative use, large scale, sensitive data, vulnerable individuals, etc.). | ☐ |
| | | | Include **specific AI risks** (automated discrimination caused by a bias in the system introduced during development, risk of producing fictitious content about a real person, risks linked to known attacks specific to AI systems, etc.) and take appropriate measures. | ☐ |

CNIL.