# Relativity®

# AI for PI: Find and Redact Personal Information with AI-Powered Workflows

# Contents

# Introduction

Meta recently faced $1.3 billion in fines imposed by European privacy regulators. While non-compliance with data privacy regulations can be costly, as demonstrated by the billions of dollars that have been paid out in the last few years, the reputational damage can have potentially more egregious consequences for the company even if its impact is harder to quantify. In this fast-evolving privacy landscape, it becomes imperative to have the technology designed to mitigate the risks associated with sensitive information, particularly PI. Now, more than ever, it's critical to leverage the power of artificial intelligence (AI) to find and redact sensitive information. For computationally intensive tasks involving voluminous amounts of data, AI consistently outperforms humans on accuracy and eliminates the risks associated with the element of human error. Using AI in the PI identification process not only improves your output, but it also ensures the sensitive information buried in your data sets remains protected. Additionally, it saves you time, money, and a significant amount of organizational risk to power your workflows with AI.

Redaction is the process of obscuring or removing sensitive information, such as text or images, from a document prior to publication or release — whether in the context of a document production for a legal matter, a Data Subject Access Request (DSAR) for data privacy regulation compliance, or even cross–border data protection. Enterprises and the counsel representing want to redact information, to protect their organization from the risk of sharing PI, PHI, or PII and violating privacy law. Organizations redact information that is personally identifiable, sensitive, or proprietary such as trade secrets or financial data.

Redaction protects an organization and the people it serves by removing sensitive information. However, as increasing amounts of data are stored and collected, and enterprises deal with new categories of data privacy regulation, manual approaches to redaction have become increasingly difficult to manage. Manual approaches may not only necessitate additional steps in document production, but they can also be marred with inconsistencies in choosing what data is redacted.

# Why is Automated Identification of PI Better?

In tandem with the explosion of data that enterprises collect, store, and process, the number and breadth of data privacy regulations that govern how enterprises need to account and share that data has proliferated. The outcome is that there are more consequences to not ensuring redaction of sensitive and personal data in documents shared with third parties, while existing manual approaches by their nature cannot scale to address growing volume.
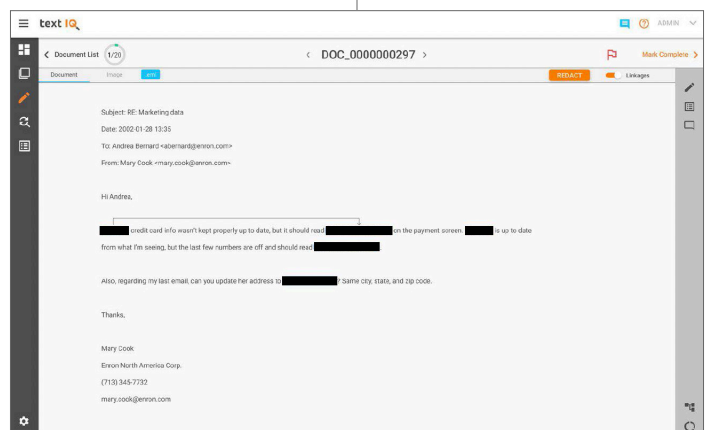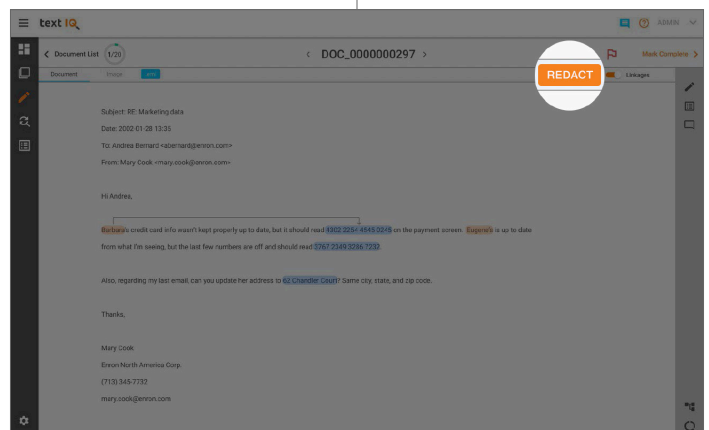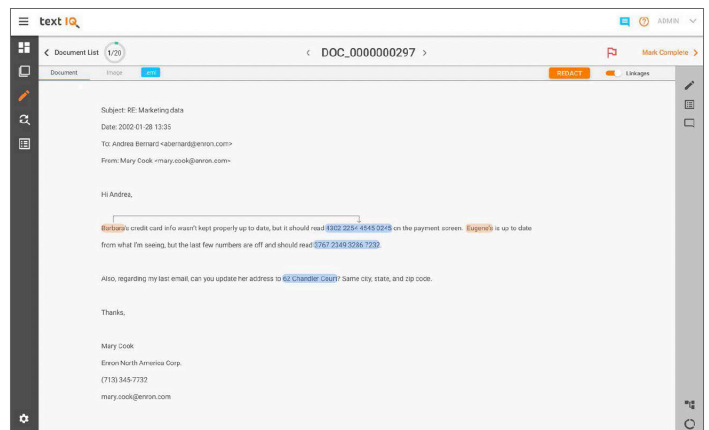
AI can uncover patterns in data that human reviewers can't on their own, allowing you to find critical PI that traditional searches miss, while also ensuring accuracy and consistency. When the stakes are high, you can't leave PI identification to chance. While each case is unique, there are a set of standard identifiers that will need to be redacted.

# Problems with Manual PI-Identification

Before the digital age, redaction meant inking over or cutting out information from a document, and — believe it or not — there are still some people who print out their documents, redact by hand, and then scan them back into their computer systems. It goes without saying that this is woefully inefficient.

More recent methods involve redaction software that provide a user interface for searching through documents and making redactions. While we can look for keywords or even black out every instance of specific words or phrases, this approach has obvious drawbacks:

Not only does it take a lot of time to redact sensitive information hiding in tens of thousands of emails or from databases in the terabyte range, it's also

easy to overlook instances of sensitive information that need redaction. For example, a search for the keyword "address" may miss the typos "adress" or "addres" made during data entry, or it could miss the phone number with a 123-456-7890 format among all the numbers with (123) 456-7890 formatting (i.e., formatting that involves parenthesis). Additionally, there is the very real risk of missing codewords, exceptions, and other irregularities.

Even if an application is programmed for auto-redaction, solving the time issue, the larger, more critical accuracy problem remains. The simple fact is that human-generated data is far too unstructured and noisy for rule-based programs to handle. This leaves organizations and firms in a bind: unwilling to rely on imprecise software to complete the job, and unable to afford the cost and risk associated with human reviewers.

## When Does PI Need to Be Redacted?

There are several instances where PI needs to be redacted:

**Litigation:** Litigation includes several steps, but there are four instances that require PI-identification: Responsiveness, Issues, PI Identification, and Privilege. The most prevalent method for PI Identification is manual, which is both costly and inefficient.

**Privacy:** Accurate PI identification is required for compliance with GDPR, CCPA, and other privacy regulations. Quickly and easily finding PI enables organizations to fulfill data subject rights and implement appropriate privacy safeguards. There are higher standards in the EU due to more stringent privacy regulations, which makes PI identification and redaction a must-have anytime they are moving data across borders.
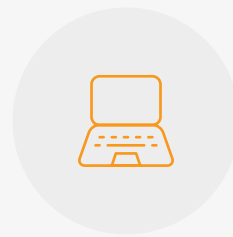
**Business Sensitive:** Proprietary information and business confidential information must be found and redacted as part of M&A, particularly for highly regulated organizations.

**Data Subject Access Requests (DSAR):** Timely retrieval of specific individual data and any identifying data that is not related to the data subject needs to be redacted.

**GOOD**
Manually redacting personal information

**BETTER**
Software that redacts based on search terms, keywords, and regular expressions

**BEST**
AI-based personal information identification and redaction

## What Type of Information Needs to be Found?

Personally identifiable information needs to be found and redacted to comply with the growing number of privacy regulations. Failure to do so can result in heavy fines. Here are some examples covered under GDPR:

**Personal Identifiers:** Name, address, SSN, phone number, birthdate, fax number, email address

**Financial Details:** Account numbers and information, credit card numbers

**Vehicle Identifiers:** VINs, license plate numbers, serial numbers

**Geographic Identifiers:** Geographic subdivisions smaller than states, initial 3 digits of zip codes

**Personal Health Information:** Names, dates (except year), age, phone numbers, fax numbers, email addresses, medical record numbers, health plan beneficiary numbers, account numbers, certificate/ license numbers, data held by a hospital or doctor

**Device Identifiers:** IP addresses, serial numbers, URLs

**Biometric Identifiers:** Including finger and voice prints
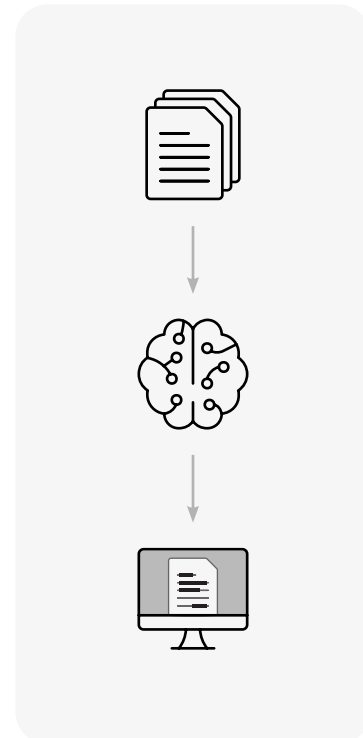
**Photographs:** Full-face photos and any comparable images

**Other:** Judiciary records, Gender. Political opinions. Identification card number. A cookie ID, Internet Protocol (IP) address, Location data (for example, the location data from a mobile phone). The advertising identifier of your phone.

# Artificial Intelligence for PI-Identification

Personal Information uses machine learning to go beyond regular expressions and to understand context and find PI that traditional approaches miss. It has already been a huge game changer for large enterprises and law firms. For example, during the document review phase of a large corporate litigation, Personal Information reduced time spent on redaction by 78 percent, which equated to an approximately $100,000 cost reduction and a six-week time reduction for the project.

How can AI achieve these breakthrough results? Unstructured data, such as emails, text documents, instant messages, photos, audio files, and other types of unlabeled data can't be accurately analyzed using rule-based applications. Personal Information utilizes a type of AI called self-supervised learning. These algorithms begin training on clean data, like the type that a search function could easily find, and then use that learning to expand into the chaotic noise that makes up most of our data sets. Instead of relying on subject matter experts to tell it the rules, self- supervised learning figures them out on its own.

Once the parameters for the types of PI that you need to identify are set, the product independently analyzes the data, finds personal information, and your team can run a quality check and quickly redact. It's quicker, easier, and more reliable than any other method.

---

"Anywhere we can reduce the burdensome process of redaction would be extremely helpful. Hearing things like auto-redaction gets me very excited."

**LEEANNE MANCARI**

Litigator, Co-Chair of eDiscovery & Information Management Platform,
DLA PIPER

---

# The Takeaway

Technology excels at automating routine, tedious tasks, freeing people up to do work that is more meaningful and valuable. Between new regulations and exploding data volumes, it's more important than ever to leverage AI across workflows. AI-based personal information identification was historically too complicated for software to handle accurately due to the inability to understand context, semantics, and other intricacies, as well as the high risk of failed regulatory compliance. We couldn't trust them with independent redaction. Fortunately, AI is now smart enough to tackle these complex tasks.

**Interested in learning more?**

Email **sales@relativity.com** with questions or to schedule a demo.

**⊞Relativity**®