# ABBYY Smart Classifier

## Technical Product Overview

## The Challenge of Classifying Unstructured Information

For storing data, searching through documents and databases by keywords or defined expressions, computers are useful tools. If a piece of text contains the search phrases or unique keywords, computer-based rules can easily be written and applied to help organise and manage content assets in databases. However, rules do not work well for unstructured content written in natural language. Documents such as articles, reports, process descriptions, product documentation, letters, and e-mails, can contain information also in free flowing text and may not necessarily use the exact phrases to make computer-based search or classification possible.

Enterprises and organisations have collected hundreds of thousands of documents which carry relevant information – but automatically identifying subject area, department, or process to which each of them belong, is still an IT challenge.

Structures and categories help the human brain to organise the world. In the same way, categorising the content of documents being stored or arriving at enterprises can greatly improve their process efficiency. In order to enable existing IT systems to classify unstructured content as reliable as structured documents, language-based technology that derives classification categories automatically from text is needed.

## What is Smart Classifier?

ABBYY Smart Classifier is a scalable server-based classification application for organising unstructured information based on statistics, morphology and semantics. Integrated via REST API this classification application becomes an intelligent component within existing IT systems, workflows, and solutions.

Smart Classifier's set-up and classification model optimisation is intuitive. The system automatically selects and optimises the classification features and parameters. Consequently, Smart Classifier delivers high-quality classification results without the need for data experts or scientific classification background. Smart Classifier is production-ready and natively supports 39 languages and a large variety of file formats such as text, office document, PDF, and image formats.

## Approaching the Unstructured Data Challenge

ABBYY Smart Classifier can "read" a selection of pre-classified documents. During the training phase the documents will be analysed to determine what document content for each particular category has in common. Once the system is trained and deployed, content assets can be sent for classification. Smart Classifier will analyse the new assets by comparing relevant features and content with the documents it was trained on.

The result is a statistical classification decision which is based on the most applicable classes that the document will belong to. The generated metadata from each document can be used to further optimise information and knowledge management systems or business processes. The scope of use cases how classification helps to automate information processes is vast and include: to determine whether a piece of text is confidential or not, add automatically generated tags to documents stored in repositories, route documents to the relevant downstream business process, workflow or department, or enhance search by restricting it to a particular class of records.

## HIGHLIGHTS & BENEFITS

Linguistically enhanced text- and semantic-based classification for managing unstructured information and documents

High quality classification results out-of-the-box

Straightforward set-up and creation of classification models with a webbased Model Editor - no "data science" expertise required

Easy set-up and training of new classification models with a web-based Model Editor

Automatic classification algorithm optimisation

Flexibility to address new and changing process needs. Fast creation of new classification "classes" using small sample training sets

Native support of many file formats such as Office, PDF and image

Embedded optical character recognition (OCR)

Ready for international projects with linguistic classification support for 39 languages

REST APIs for easy integration

Scalable processing backend

## Getting started

ABBYY Smart Classifier is flexible and capable of processing content from different subject areas, domains and departments.

The system learns automatically based on the training documents. Every organisation has experts with knowledge of the internal documents that are generated or received within that organisation. These people are perfectly suited to train classes in Smart Classifier. No in-depth scientific knowledge is required to generate high classification quality. The system extracts the text from submitted documents and analyses it using the morphology of the language, semantic and textual features, statistics and machine learning.

The principle behind machine learning is to identify and use the most relevant features automatically from a set of training documents.

## Preparation: Category Definition

Content and process experts within the organisation are perfectly suited to define and train the necessary categories (or classes) in Smart Classifier and select the training documents - prototypical sample documents that represent the categories.
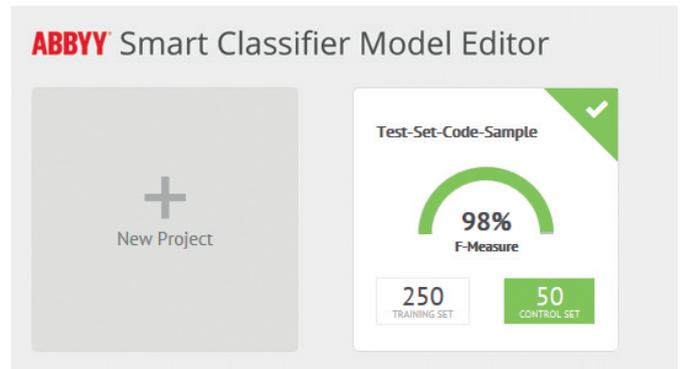
With Smart Classifier, there is no need to use thousands of documents to train a category. The minimum technical number is ten documents for a class, while a set of 100+ documents is recommended to generate reliable statistics. Smart Classifier makes collecting training documents easy since it can natively process a large variety of file formats present in enterprise reality such as plain text, office formats, HTML, etc. Thanks to built-in ABBYY OCR and conversion capabilities, the application can also process PDFs, faxes or scanned documents reliably. The training documents for the different classes should be set up for individual languages and arranged in a folder structure.

Now the training document collection from each category should be divided into a training set and a control set. Once compressed as zip-archives they can be uploaded via the intuitive Model Editor interface to train the system.
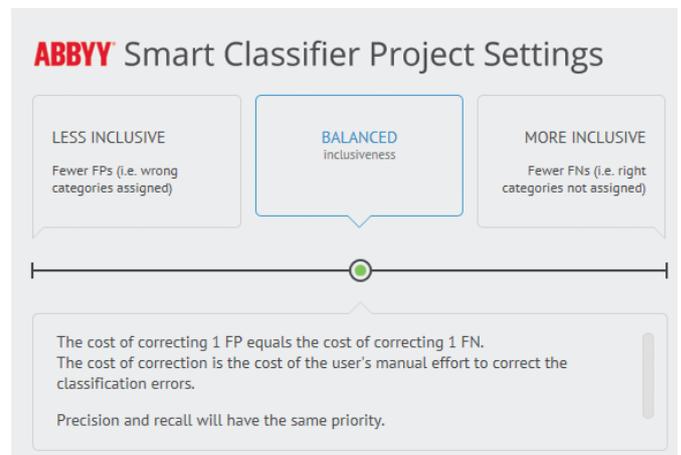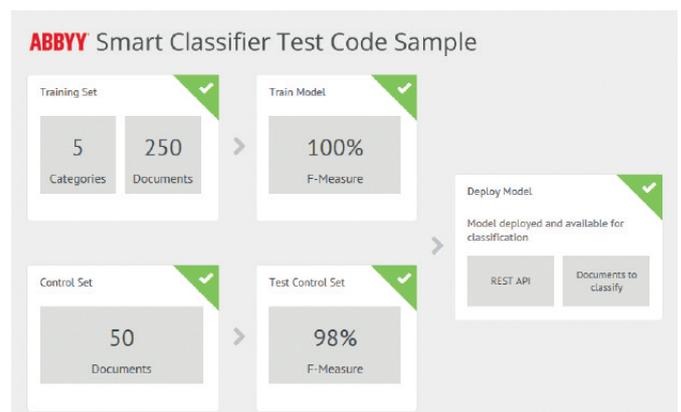


## System Training: Model Editor

Smart Classifier's Model Editor offers an intuitive and easy-to-use web interface. Classification projects, document training sets and models are created and managed remotely via a browser.



New classification projects can easily be set up by specifying the language for a model language and the inclusiveness of classification. With these settings, the sensitivity of the model is controlled. Depending on the business need, it is possible to adjust the model if the classification should "focus" on a high "true positive" rate or if the model can be more "open" – with a statistical risk of putting too many documents in a particular category (= precision/recall balance).



The Model Editor provides a status overview for each classification project and gives access to the different components of the workflow including project settings, training sets, control set documents and quality evaluation for each of the classified assets. Unknown or wrongly classified documents can be uploaded to evaluate/debug the classification results and words, even after the model is deployed.

## Quality Evaluation

Smart Classifier makes classification easy because classification model tuning on an algorithmic level is not required. However, since the underlying core technologies are very complex, it is important to review and determine whether the training process was successful and that the classification results in a control set meet expectations before the model is deployed. Performance metrics such as f-measure, precision/recall, true/false positives, and more help to evaluate the classification model quality.

The classification results of the documents used in the training and control sets can be checked and re-assigned if the category does not match.

**ABBYY Smart Classifier Document Classification**

| File name | Best Category | | 2nd Category | | 3rd Category |
|---|---|---|---|---|---|
| 53068.txt | sci.space | 11% | sci.electronics | 9% | rec.autos |
| 54054.txt | rec.sport.hockey | 99% | sci.med | 0% | rec.autos |
| 54060 D612.t... | rec.sport.hockey | 95% | sci.med | 24% | rec.autos |
| 54154.txt | sci.electronics | 73% | rec.autos | 16% | sci.med |
| 54180.txt | sci.electronics | 99% | sci.space | 58% | rec.autos |

**ABBYY Smart Classifier Testing Report**

The Model Editor also provides instant visibility of each document within a classification project. Source text and key words picked by the algorithms can be analysed and checked. Terms that should be ignored during classification can be added to a stop word list.

**ABBYY Smart Classifier Control Set**

## Architecture & Integration

### Simple Setup & Configuration

Smart Classifier is built on the ABBYY Compreno processing backend. All necessary components, including Microsoft Internet Information Server (IIS), can be automatically configured during the installation process. IIS hosts the web based Model Editor and allows integration via REST API.

### Backend Architecture & Scalability

Smart Classifier is based on a scalable backend, capable of processing large amounts of files. The system will extract plain text out of content assets submitted for classification; for images or PDFs, optical character recognition can be automatically applied. For a high throughput, Smart Classifier can be scaled up with additional processing stations running OCR processes.

### Integration via REST API

Smart Classifier can be easily integrated into document workflows, archiving, records management, e-mail management, or data migration systems via its RESTful API.

Once the system is set up and a classification model is published for operation, incoming classification tasks will be accepted. Depending on the amount and complexity of tasks they can be submitted either in synchronous and asynchronous modes.
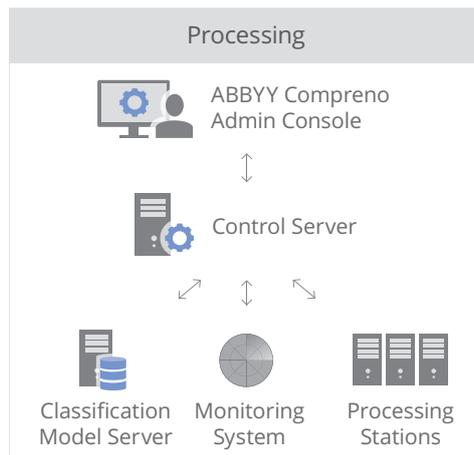
### Classification Results

Smart Classifier can return the classification results in JSON or RDF/XML. The results include information such as name of the classification model, categories with their probabilities, confidentiality flags, feature/word lists, access to the raw text** or error messages.

| Existing IT | ABBYY Smart Classifier Architecture | |
|---|---|---|
| Email | **Processing** | **Setup & Training** |
| Workflow | ABBYY Compreno Admin Console | |
| CRM | Documents → | |
| DMS | REST API | ABBYY Smart Classifier Model Editor |
| Database | ← Classification Results | |
| Archive | Control Server | |
| | Classification Model Server · Monitoring System · Processing Stations | |

# ABBYY Compreno, Smart Classifier & InfoExtractor

ABBYY Compreno natural language processing (NLP) technology enables businesses to understand unstructured information. Its intelligent technology understands the meaning of words and defines relationships between them. Based on these relationships, it creates semantic representations that enables text to be analysed by computers for accurate information extraction, classification and intelligent search. This language-based approach for analysing unstructured information creates new opportunities to action information and to optimise critical business processes where rule-based approaches fail.



**ABBYY Smart Classifier** organises unstructured information in existing or incoming documents. It extracts and analyses text with linguistic and statistical methods and derives the highest probability for the best matching classes.

**ABBYY InfoExtractor** uses a semantic based approach to extract relevant information. The system can identify not only entities and facts, but also the relationships between them. The technology is currently available for English and Russian language.

# Licensing

Smart Classifier is available for testing via time and volume limited trial licences.

ABBYY offers 3 different licensing models for Smart Classifier:
- perpetual licences with software maintenance;
- yearly subscriptions and;
- OEM licensing models for software vendors.

The standard license model is based on a renewable peak volume. The back-end can be scaled up as it is needed; the number of processing stations/CPU cores is not limited.