

ABBYY FineReader : présentation de quelques fonctions d'amélioration des performances

Par ABBYY

Date de publication : 28 janvier 2020

TOUT PUBLIC

Le kit de développement logiciel ABBYY FineReader Engine permet aux développeurs et aux fournisseurs de logiciels de créer des applications qui peuvent extraire automatiquement des données de texte dans les documents, des captures d'écran et des interfaces machine (formulaires, etc.).

Le but de ce tutoriel est de vous présenter plusieurs techniques pour améliorer la qualité de la reconnaissance et la vitesse de traitement. **Commentez**

I - Introduction.....	3
II - Facteurs influençant la précision et la vitesse de traitement de l'OCR.....	3
II-A - Type et qualité d'image.....	3
II-A-1 - Comment obtenir des images de bonne qualité.....	4
II-A-1-a - Conseils de numérisation de documents.....	4
II-A-1-b - Conseils pour prendre des photos.....	4
II-A-2 - Amélioration de la qualité d'image avec FineReader Engine.....	5
II-A-3 - Profils de traitement prédéfinis dans FineReader Engine.....	5
II-B - Modes de reconnaissance.....	6
II-C - Langues des documents.....	7
II-D - Étapes de traitement de documents.....	7
II-D-1 - Importation de l'image.....	7
II-D-2 - Prétraitement de l'image.....	8
II-D-3 - Analyse de documents.....	8
II-E - Traitement parallèle avec plusieurs cœurs.....	8
III - FineReader Engine - Résultats des tests de vitesse.....	9
IV - Ressources système.....	9
V - Comment augmenter la vitesse de traitement globale dans FineReader Engine.....	10
VI - Comment améliorer la qualité de reconnaissance de texte dans FineReader Engine.....	10
VII - Sources d'informations complémentaires.....	11
VIII - Remerciements Developpez.com.....	11

I - Introduction

L'intégration de la technologie de reconnaissance optique de caractères (OCR) élargira les fonctionnalités de votre application de manière efficace. L'excellente performance du composant OCR est l'un des facteurs majeurs du taux élevé de satisfaction des utilisateurs.

Ce tutoriel fournit des informations sur les facteurs généraux de performance de l'OCR et leurs possibilités d'optimisation dans le Kit de développement logiciel ABBYY FineReader Engine. En exploitant ses capacités et options avancées, il est possible d'améliorer encore la performance déjà élevée de l'OCR en vue d'une expérience client optimale.

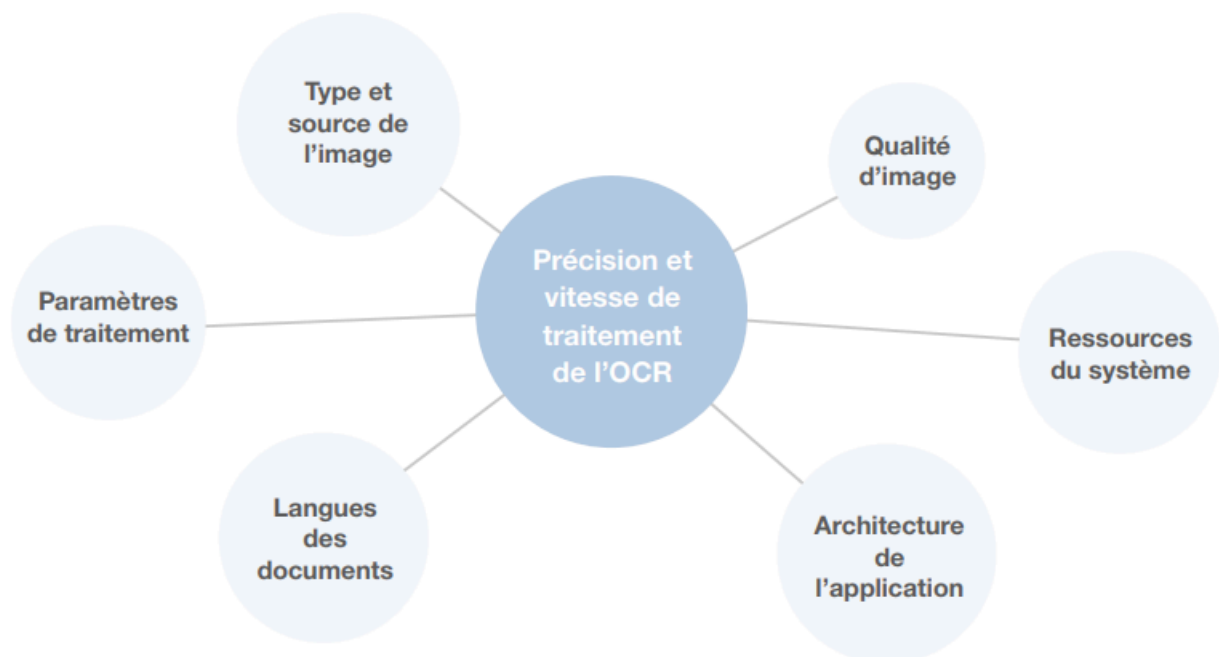
Deux paramètres essentiels sont à prendre en considération pour mesurer la performance de l'OCR :

- **la précision de reconnaissance ;**
- **la vitesse de traitement.**

II - Facteurs influençant la précision et la vitesse de traitement de l'OCR

Il est possible d'améliorer la vitesse et la précision de la reconnaissance de manière considérable en utilisant les bons paramètres dans ABBYY FineReader Engine.

Les facteurs clés qui influencent la précision et la vitesse de l'OCR sont présentés comme suit.



II-A - Type et qualité d'image

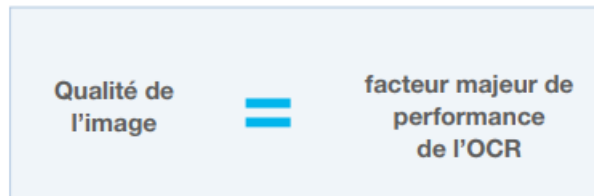
Les images peuvent provenir de différentes sources :

- fichiers PDF créés par numérisation ;
- copies d'écran d'ordinateur ou de tablette ;
- fichiers images créés par scanner ;
- serveurs de fax ;
- appareils photo numériques ou smartphones...

Différentes sources produiront différents types d'images avec des niveaux de qualité variables.

Par exemple, l'utilisation des mauvais paramètres d'un scanner peut provoquer du « bruit » sur l'image, comme des points noirs ou des taches à des emplacements aléatoires, des lettres floues et de forme irrégulière, des lignes déformées et/ou des bordures de tableau décalées. En termes d'OCR, il s'agit d'« images de mauvaise qualité ».

Le traitement d'images de mauvaise qualité nécessite une grande puissance informatique, augmente le temps de traitement global et compromet les résultats de la reconnaissance.



Par ailleurs, le traitement d'« images de bonne qualité » sans distorsions réduit le temps de traitement. De plus, la lecture d'images de bonne qualité donne des résultats dont la précision est plus élevée.

Il est par conséquent recommandé d'utiliser des images de bonne qualité pour le traitement OCR.



Augmentez la vitesse et la précision de l'OCR en améliorant la qualité de l'image.

II-A-1 - Comment obtenir des images de bonne qualité

II-A-1-a - Conseils de numérisation de documents

Taille de police

Les documents imprimés dans des polices très petites doivent être numérisés avec des résolutions plus élevées.

Utilisez la résolution suivante pour la numérisation :

- 300 dpi pour les textes standard (imprimés en polices de taille 10 ou supérieure) ;
- 400-600 dpi pour les autres textes (taille de police 9 ou inférieure).

Qualité d'impression

Un document de mauvaise qualité, tel que des journaux ou des livres anciens, doit être numérisé en mode « Niveaux de gris ».

Ce mode conserve davantage d'informations sur les lettres du texte numérisé.

II-A-1-b - Conseils pour prendre des photos

Luminosité correcte

- Assurez-vous que la luminosité est répartie régulièrement sur l'ensemble de la page et qu'il n'y a pas de zones sombres ou d'ombres.

- Si possible, utilisez un trépied. Positionnez l'objectif parallèlement à la surface du document et dirigez-le vers le centre du texte.
- Désactivez le flash pour éviter les reflets et les ombres marquées sur la page.
- Si l'appareil photo dispose d'une option de « balance des blancs », utilisez une feuille de papier blanche pour régler la balance des blancs. Sinon, sélectionnez le mode de balance des blancs qui convient le mieux aux conditions de luminosité ambiante.

Pas suffisamment de lumière

- Sélectionnez une valeur d'ouverture supérieure.
- Sélectionnez une valeur ISO de sensibilité supérieure.
- Utilisez la mise au point manuelle si l'appareil photo n'est pas capable de verrouiller automatiquement la mise au point.

Pour plus d'informations, consultez la rubrique d'aide du développeur de FineReader Engine :



Guided Tour → Best Practices → Source Image Recommendations

II-A-2 - Amélioration de la qualité d'image avec FineReader Engine

S'il n'est pas possible d'agir à l'avance sur la qualité de l'image, il est recommandé de l'améliorer avant la phase de reconnaissance. Dans FineReader Engine, différentes fonctions de prétraitement de l'image sont disponibles.

Fonctions de base :

- mise à l'échelle de l'image ;
- recadrage de l'image ;
- découpage de l'image ;
- redressement des lignes ;
- reflet et inversion ;
- suppression du bruit ;
- amélioration du contraste local ;
- correction des distorsions géométriques ;
- binarisation adaptable.

Fonctions avancées :

- technologie appareil photo OCR d'ABBYY ;
- division automatique des pages doubles ;
- suppression de cachets et notes manuscrites ;
- redressement automatique de l'image ;
- détection automatique de l'orientation des pages et rotation ;
- élimination des taches de l'image.

II-A-3 - Profils de traitement prédéfinis dans FineReader Engine

La précision de reconnaissance et la vitesse du traitement peuvent être optimisées en affinant les paramètres individuels. Cependant, on peut également appliquer un des profils de traitement prédéfinis, qui sont disponibles pour les scénarios d'utilisation courants. Les configurations fournies dans ces profils sont adaptées à la plupart des situations.

La plupart des profils sont proposés sous deux formes :

- avec des paramètres optimisés pour la meilleure précision de reconnaissance ;
- avec des paramètres optimisés pour la vitesse de traitement la plus élevée.

Si nécessaire, il est toutefois possible d'affiner le réglage des profils en configurant d'autres paramètres utiles via l'API.

Scénario	Nom du profil
Conversion de documents pour archivage Un document nécessite une vitesse de traitement élevée, une bonne qualité visuelle et une taille réduite du fichier PDF final. Une précision de reconnaissance supérieure n'est pas un paramètre essentiel, un niveau satisfaisant est suffisant.	DocumentArchiving_Accuracy DocumentArchiving_Speed BookArchiving_Accuracy BookArchiving_Speed
Conversion de document en vue d'une réutilisation du contenu Dans ce cas, la précision de reconnaissance et la conservation de la structure du document sont les critères les plus souhaités. Comme le document final doit être exempt d'erreur, un faible taux de reconnaissance ou une reconstruction erronée de la mise en page entraînerait un travail supplémentaire pour les opérateurs.	DocumentConversion_Accuracy DocumentConversion_Speed
Extraction d'un texte sur une image comprenant de petites zones de texte de mauvaise qualité. Par la suite, les informations de ce texte doivent pouvoir être recherchées, extraites pour un autre traitement ou utilisées pour classer le document.	TextExtraction_Accuracy TextExtraction_Speed
Capture de données à partir de zones définies précisément dans la page. Reconnaissance de codes-barres.	FieldLevelRecognition BarcodeRecognition_Accuracy BarcodeRecognition_Speed
Création de fichiers PDF ultracompressés, enregistrés sous forme d'images. Reconnaissance de cartes de visite.	HighCompressedImageOnlyPdf BusinessCardsProcessing
Reconnaissance de dessins techniques, généralement de taille importante et qui comportent des diagrammes complexes ainsi que des orientations de texte différentes.	EngineeringDrawingsProcessing

II-B - Modes de reconnaissance

Pour agir sur la performance de l'OCR, il est également possible d'utiliser des modes de reconnaissance conçus pour des scénarios spécifiques. FineReader Engine fournit les modes de reconnaissance prédéfinis suivants :

- **Mode normal** : via ce mode, vous obtiendrez la plus haute précision de reconnaissance. Ce mode est fortement recommandé lorsque la reconnaissance d'un contenu est destinée à être réutilisée dans d'autres applications ou tâches pour lesquelles la précision est d'une importance primordiale ;
- **Mode équilibré** : ce mode fournit des valeurs intermédiaires de précision de reconnaissance et de vitesse. En général, ce mode de reconnaissance offre une vitesse supérieure au mode de reconnaissance « normal » tout en atteignant quasiment le même niveau de précision ;
- **Mode rapide** : l'utilisation de ce mode augmente la vitesse de traitement jusqu'à 200-250 %. Ce mode est recommandé lorsque la vitesse de traitement est une priorité, dans le cadre par exemple du traitement d'un document très volumineux pour archivage, de systèmes de gestion de contenus et documentaire.

II-C - Langues des documents

ABBYY FineReader Engine est capable de reconnaître aussi bien des documents monolingues que multilingues (c'est-à-dire rédigés en plusieurs langues). Il est très important de définir la langue de reconnaissance adéquate, faute de quoi le traitement du document pourrait être considérablement ralenti et la qualité de reconnaissance diminuée.

Si la langue de reconnaissance ne peut être définie à l'avance, il est possible d'utiliser la détection automatique de la langue.

Cependant, la présélection d'un nombre important de langues de reconnaissance réduira la vitesse de traitement.

Il est par conséquent déconseillé de définir plus de cinq langues de reconnaissance.

Pour accroître encore la précision de reconnaissance, FineReader Engine fournit un dictionnaire et une assistance morphologique pour de nombreuses langues. Lorsque le traitement de documents comporte des termes spécifiques à un sujet ou des « structures » telles que des codes produits, des numéros de téléphone ou des numéros de passeport, des dictionnaires créés sur mesure peuvent être importés pour garantir une qualité de reconnaissance élevée.

II-D - Étapes de traitement de documents

En termes d'OCR, le traitement de documents est un processus en plusieurs étapes. En fonction d'un scénario particulier, des fonctions et paramètres différents peuvent être appliqués à chaque stade de l'OCR. Dans FineReader Engine, le processus comprend les étapes suivantes :

- **Importation d'image** : importation du document image à FineReader Engine pour traitement ;
- **Prétraitement de l'image** : amélioration de la qualité de l'image avant le traitement OCR ;
- **Analyse du document** : détection d'objets individuels, tels que du texte, des images ou des codes-barres ;
- **Reconnaissance** : la « lecture » extrait les textes tapés à la machine et manuscrits ainsi que les valeurs des marques optiques et des codes-barres ;
- **Synthèse** : reconstruction de la mise en page originale du document, y compris l'ordre et la numérotation des pages, le flux logique du texte, les images, etc.
- **Export** : exportation finale des résultats du traitement dans les formats requis.

Les différentes méthodes et paramètres utilisés pour chaque scénario de traitement influenceront considérablement la vitesse globale du traitement. Il est possible d'accélérer l'ensemble du processus d'OCR en supprimant les étapes inutiles. Par exemple, lorsqu'on extrait des données à partir de zones prédéfinies d'un document, l'analyse du document n'est pas nécessaire. Lors de l'exportation de documents aux formats TXT ou PDF « image uniquement », l'étape de synthèse peut être ignorée.

II-D-1 - Importation de l'image

Il est possible d'envoyer directement des documents dans FineReader Engine à partir d'un scanner ou importés à partir d'un système de stockage ou d'un flux de mémoire.

Différentes méthodes sont nécessaires pour obtenir des images de différentes sources et influencer la vitesse de reconnaissance. L'importation d'images à partir de la mémoire est généralement plus rapide que l'ouverture des images à partir d'un fichier de stockage.

II-D-2 - Prétraitement de l'image

En général, le processus d'OCR est plus rapide avec des images de bonne qualité. Il est recommandé d'optimiser l'étape de prétraitement de l'image en conséquence, ce qui permettra de gagner du temps lors de l'étape de traitement.

Les images peuvent être de différents formats et de qualité inégale. Les images de bonne qualité, telles que les fichiers PDF créés numériquement, ne demandent pas, en principe, beaucoup de travail préalable. Pour les images de mauvaise qualité, comme les documents scannés avec des paramètres de scan non adaptés ou des livres anciens, il est nécessaire d'appliquer les fonctions avancées de prétraitement d'image pour améliorer les résultats de la reconnaissance.

Pour le prétraitement de photos numériques, la technologie spéciale Camera OCR™ d'ABBYY est utilisée. Ici, les algorithmes sont optimisés spécifiquement pour améliorer les photos.

L'utilisation de différentes fonctions de prétraitement agira spécifiquement sur la vitesse de traitement.

II-D-3 - Analyse de documents

Les brochures ou les journaux comportent souvent du texte présenté en colonnes, tableaux, diagrammes ou images.

Les dessins techniques peuvent se présenter sous forme de grands documents avec des diagrammes d'ingénierie complexes et des textes avec différentes orientations.

Pour les documents comportant ce type de mise en page complexe, l'étape d'analyse de documents nécessitera plus de temps. En revanche, l'analyse de documents à la mise en page simple, comme les lettres ou les contrats, est très rapide.

II-E - Traitement parallèle avec plusieurs cœurs

FineReader Engine peut être utilisé pour concevoir des applications, indépendamment de leur taille et de leur complexité, qu'il s'agisse d'un poste de travail client, d'une solution basée sur des serveurs ou d'un projet conséquent intégrant plusieurs millions de pages. En termes de performance OCR, il est souvent judicieux d'utiliser des systèmes multiprocesseurs ou multicœurs pour augmenter la vitesse de traitement.

Le support multicœur intégré dans FineReader Engine permet de développer le processus d'OCR selon différentes méthodes :

- utilisation d'une instance Engine simple ;
- chargement d'instances Engine multiples.

Il existe différentes méthodes pour traiter les documents :

- **Traitement de gros documents multipages**

L'objet « `FRDocument` » est le mieux adapté au traitement de gros documents comportant de nombreuses pages.

Dans ce cas, les pages d'un document multipage sont traitées simultanément sur les processeurs disponibles.

À la fin, les résultats sont assemblés en un seul document multipage. C'est le multitraitement le plus facile à coder. Le nombre de traitements nécessaires est détecté automatiquement, en fonction de facteurs tels que le nombre de processeurs physiques ou logiques disponibles, le nombre de processeurs disponibles dans le cadre de la licence et le nombre de pages du document. Si nécessaire, le développeur peut facilement changer les paramètres de multitraitement et adapter le nombre de traitements à effectuer.

- **Traitement de plusieurs documents d'une seule page**

Pour traiter simultanément plusieurs documents d'une seule page, provenant de la même source, par exemple, un scanner, il est recommandé d'utiliser la méthode : « BatchProcessor ». Cette méthode est la plus efficace en termes de rapidité, si l'export du document en format PDF n'est pas nécessaire, comme dans les scénarios de capture de données avec un format de sortie personnalisé.

Pour traiter simultanément plusieurs documents d'une seule page, il est recommandé d'utiliser plusieurs instances Engine. Cette méthode est aussi la mieux adaptée à des scénarios de services Web, lorsque le document entrant doit être traité directement après réception. Dans ce cas, le document est transféré à une instance FineReader Engine du pool et traité immédiatement.

III - FineReader Engine - Résultats des tests de vitesse

Le tableau présente les résultats des tests internes de performance. Veuillez noter que les résultats des tests dépendent de nombreux facteurs tels que la qualité de l'image, les langues de reconnaissance utilisées et d'autres facteurs.

	Documents d'une page (Pages/ minute)	Un document multipage (Pages/ minute)	Scénario de capture de données sans export (Pages/ minute)
Traitement séquentiel	60	51	87
Traitement parallèle avec FRDocument	41	117	57
Traitement parallèle avec FRDocument tout en conservant les données en mémoire pendant le traitement	55	141	82
Traitement parallèle avec Batch Processor	99	115	294
Traitement parallèle avec plusieurs Engines	165	10	102

Informations techniques sur les tests

Intel®Core™ i5-4440 (3,10 GHz, 4 cœurs physiques), 8 Go RAM, 4 traitements simultanés.



La performance a été testée sur 300 documents en anglais, avec le profil prédéfini « DocumentArchiving_Speed ». Dans ces scénarios « Documents d'une page » et « Un document multipage », les documents ont été exportés au format PDF.

IV - Ressources système

Pendant le processus OCR, toute une gamme d'algorithmes sont appliqués. Ils dépendent de la qualité de l'image, des langues du document, de la complexité de la mise en page et du nombre de pages du document. Selon le cas,

de tels algorithmes peuvent nécessiter des ressources mémoire supérieures. Il est recommandé de configurer le système conformément aux exigences de mémoire spécifiées pour optimiser la vitesse de traitement en allouant la mémoire système adéquate.

Mémoire requise

- **Traitement de documents d'une page :**
 - minimum 400 Mo RAM ;
 - 1 Go RAM recommandé.
- **Traitement de documents multipages :**
 - minimum 1 Go RAM ;
 - 1,5 Go RAM recommandé.
- **Traitement parallèle :**
 - 350 Mo RAM x nombre de cœurs de processeur ;
 - + 450 Mo RAM supplémentaires.
- **Traitement parallèle de documents en arabe, chinois, japonais ou coréen :**
 - 850 Mo RAM x nombre de cœurs de processeur ;
 - + 750 Mo RAM.

V - Comment augmenter la vitesse de traitement globale dans FineReader Engine

Il existe plusieurs possibilités pour améliorer la performance de votre système.

- Affinez les paramètres de prétraitement de l'image pour fournir la meilleure qualité de document pour l'étape de traitement.
- Pendant la phase de traitement, utilisez un des profils de traitement prédéfinis optimisés au niveau de la vitesse et le mode de reconnaissance adéquat – équilibré ou rapide.
- Définissez les langues de reconnaissance adéquates. Une mauvaise définition de la langue peut ralentir considérablement le traitement du document. Plus il y a de langues de reconnaissance sélectionnées, plus la vitesse de traitement est lente.
- Utilisez l'objet approprié (FRDocument ou BatchProcessor) et activez le traitement parallèle.
- Définissez les paramètres appropriés d'analyse et de reconnaissance. Par exemple, désactivez la détection de tableau et la correction de l'orientation de page si les images ne comportent pas de tableau et si l'orientation des pages est correcte.
- Ignorez l'étape de synthèse si les documents traités sont à exporter aux formats TXT ou PDF « image uniquement ».
- Utilisez le profil Fast PDF Export pour exporter des documents au format PDF.
- Utilisez l'objet spécial (ExportFileWriter), qui est conçu pour l'export de très gros documents multipages au format PDF.

Pour plus d'informations, consultez la rubrique d'aide du développeur FineReader Engine :

Guided Tour → Best Practices → Increasing Processing Speed

VI - Comment améliorer la qualité de reconnaissance de texte dans FineReader Engine

FineReader Engine propose une qualité de reconnaissance élevée. La qualité de la reconnaissance dépendra toujours de facteurs tels que la qualité de l'image, la langue et autres facteurs. Il existe toutefois plusieurs moyens d'améliorer la qualité de la reconnaissance.

- Définissez correctement le type de texte.
- Définissez les langues de reconnaissance adéquates.
- Définissez les langues rares et les dictionnaires sur mesure pour la reconnaissance de caractères spéciaux ou les documents avec une terminologie particulière, par exemple, des textes juridiques ou médicaux.
- Scindez les pages en vis-à-vis des livres scannés en deux images distinctes.
- Appliquez la technologie spéciale Camera OCR pour le traitement de photos numériques.
- Corrigez la résolution de l'image, si elle diffère considérablement de la résolution recommandée.

Pour plus d'informations, consultez la rubrique d'aide du développeur FineReader Engine :

[Guided Tour](#) → [Best Practices](#) → [Improving Recognition Quality](#)

Comme vous le savez, le processus OCR peut être une tâche très complexe du traitement de documents. En fonction du scénario de traitement de chaque document, les résultats de performance de l'OCR peuvent varier significativement. Les conseils d'optimisation de la précision de reconnaissance et de vitesse de traitement dans ABBYY FineReader Engine devraient vous aider à obtenir la meilleure performance professionnelle.

VII - Sources d'informations complémentaires

Pour en savoir plus quant aux différents aspects d'optimisation de performance OCR et sur les kits de développement logiciel ABBYY FineReader Engine, veuillez utiliser les sources d'information suivantes :

- portail technologie ABBYY : <https://abbyy.technology/> ;
- forum kits de développement logiciel OCR ABBYY : <https://forum.ocrsdk.com/> ;
- le fichier d'aide fourni à la distribution d'ABBYY FineReader Engine ;
- pages produit ABBYY FineReader Engine sur www.abbyy.com/ocr-sdk.

VIII - Remerciements Developpez.com

Developpez.com remercie ABBYY pour l'autorisation de publication de ce tutoriel. Tous les remerciements aussi à **Guillaume SIGUI** pour la mise au gabarit et **Claude Leloup** pour la relecture orthographique.