

# Apprentissage d’observateurs idéaux pour l’évaluation de la qualité d’images en Tomographie par Rayons X en vues éparses

Romain VO<sup>1,2</sup> Antoine OLISLAEGERS<sup>2</sup>

<sup>1</sup>CNRS, ENS de Lyon, Laboratoire de Physique, Lyon, France

<sup>2</sup>Université Paris-Saclay, CEA, List, F-91190, Palaiseau, France

**Résumé** – Dans cet article, nous tentons de répondre à une question primordiale en imagerie : *Comment évaluer la qualité d’une image ?* En nous concentrant sur des applications de tomographie par rayons X, nous présentons un cas limite d’évaluation à l’aide du PSNR, qui ne parvient pas à capturer la qualité visuelle d’une image (voir Fig. 1). Nous nous appuyons sur le concept d’*observateurs idéaux* pour proposer une nouvelle méthode d’évaluation de la qualité d’une image, et surtout **une métrique** capable de classer différents algorithmes de restauration en termes de *qualité* des images produites. En effet, en formulant le problème comme un *test d’hypothèse*, l’approche par *observateur idéal* offre une manière intuitive de quantifier la qualité d’une image.

**Abstract** – In this work, we try to find an answer to a longstanding question in computational imaging: *How to assess the quality of an image ?* Focusing on X-ray computed tomography (CT), we present a failure case of PSNR, which fails to capture the visual quality of a reconstruction (see Fig. 1). We build on the concept of *ideal observers* to propose a new way to assess the quality of an image, and most importantly **a metric** that is able to rank different image restoration algorithms in terms of the *quality* of the images they produce. The *ideal observer* framework provides an intuitive way to quantify the quality of an image.

## 1 Introduction

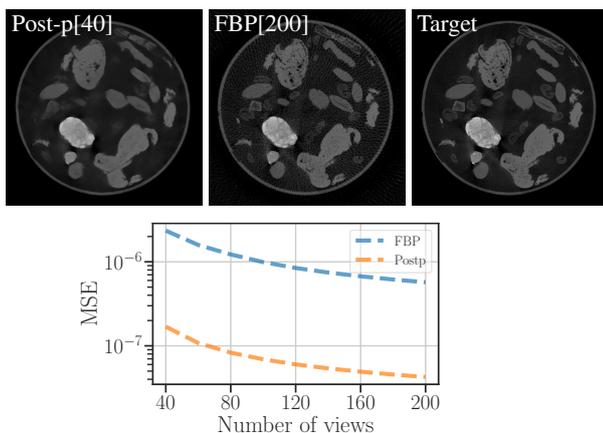


Figure 1 – Sur le jeu de données **2DeteCT** [7]: évolution de la MSE des reconstructions FBP (*courbe bleue*) et des reconstructions post-traitées par un réseau de neurones profond (*courbe orange*) en fonction du nombre de vues. De gauche à droite: une reconstruction FBP à 40 vues post-traitée par un réseau de neurones profond, une reconstruction FBP à 200 vues, et l’image cible.

L’évaluation de la qualité d’image (*Image Quality Assessment – IQA*) est une étape cruciale dans le développement de systèmes d’imagerie tels que la tomographie par rayons X (CT). Un processus d’évaluation de la qualité bien défini aide les utilisateurs à établir des normes, par la suite garantissant que les images produites sont de qualité suffisante pour la tâche à accomplir, qu’il s’agisse de diagnostic, de surveillance ou d’inspection industrielle. Cela permet également de comprendre les limites d’un système, d’identifier ses modes de défaillance et ses biais – autant d’étapes nécessaires à son amélioration.

Dans ce contexte, les métriques standards d’évaluation à partir d’images de référence (*full reference IQA – FR-IQA*) telles que le maximum du rapport signal sur bruit (PSNR) et l’indice de similarité structurelle (SSIM) sont largement utilisées pour évaluer la qualité des images [14]. Cependant, ces métriques se retrouvent dans de nombreux cas inadéquates. Le PSNR a par exemple tendance à privilégier les images lisses au détriment des images riches en détail [2].

La Fig. 1 illustre ce problème: on y observe que l’erreur quadratique (*Mean Squared Error – MSE*) d’une reconstruction à 40 vues post-traitée par un réseau de neurones profond (*post-processing*) est inférieure à celle d’une reconstruction par rétro-projection filtrée (*Filtered Back Projection – FBP*) à 200 vues, alors que cette dernière est visuellement de meilleure qualité. Ce constat est critique car il peut mener à la sélection d’algorithmes sous-optimaux. En effet, une grande partie de la recherche sur le développement de nouveaux algorithmes de reconstruction se concentre sur l’amélioration de ce type de métrique, qui ne sont pas toujours corrélées avec la réussite d’une tâche en aval sur ces mêmes images [2].

Dans cet article, nous explorons l’utilisation du critère de Bayes et plus spécifiquement des observateurs idéaux (*Ideal Observers – IO*) [15, 1, 11] comme cadre général pour évaluer la qualité d’une image. L’évaluation par observateur part de l’idée largement admise selon laquelle un système d’imagerie doit être optimisé pour une tâche spécifique [1, 2, 6], *e.g.* en imagerie médicale, cela signifie optimiser la capacité du médecin à détecter une lésion.

Pour une tâche de détection binaire, ce processus peut être formulé comme un test d’hypothèse: un observateur est présenté avec l’hypothèse nulle ( $H_0$ ), *i.e.* l’image ne contient pas de signal, et l’hypothèse alternative ( $H_1$ ), *i.e.* l’image contient un signal. Dans ce contexte, l’observateur idéal [15] fournit une formule pour calculer la statistique de test optimale et fixe une limite supérieure sur la performance de tout

observateur, humain ou algorithme, effectuant la tâche. Il est optimal car l'observateur idéal minimise la probabilité d'erreur moyenne, *i.e.* la probabilité de faire une mauvaise décision, étant donné les données [15]. Malheureusement, l'observateur idéal est analytiquement intractable et n'est utilisé que dans des modèles relativement simples où l'image et le signal à détecter suivent des distributions statistiques simples [9]. Une alternative coûteuse consiste à calculer une approximation linéaire tel que l'observateur d'*Hotelling* [3].

Des méthodes récentes explorent la possibilité d'apprendre l'observateur idéal à partir de données [8, 17, 16]. Les résultats obtenus sur des ensembles simples, avec un observateur idéal analytique, suggèrent que l'apprentissage de l'observateur idéal à partir de données est prometteur et peut être utilisé pour évaluer la qualité des images. Cependant, peu de travaux ont tenté d'utiliser des observateurs appris pour classer des algorithmes d'imagerie [16, 10].

**Contributions** Ce travail vise à fournir un aperçu clair de l'évaluation par observateur. Nous nous concentrons sur la tomographie par rayons X en vues éparses et comparons deux algorithmes de reconstruction: la rétro-projection filtrée (FBP) et une méthode de post-traitement [4] par apprentissage. Nous faisons varier le nombre de vues pour simuler différentes conditions d'imagerie et évaluons comment le post-traitement affecte la performance mesurée par un observateur appris.

## 2 Observateurs

Une tâche de détection binaire peut être formulée comme un *test d'hypothèse* : étant donné une image  $\mathbf{x} \in \mathbb{R}^n$ , un *observateur* doit choisir entre deux hypothèses: l'*hypothèse nulle* ( $H_0$ ), l'image  $\mathbf{x}$  ne contient aucun signal d'intérêt, et l'*hypothèse alternative* ( $H_1$ ), l'image contient un signal aléatoire  $\mathbf{s} \in \mathbb{R}^n$  :

$$\begin{cases} H_0 : \mathbf{x} = \mathbf{x}_b, \\ H_1 : \mathbf{x} = \mathbf{x}_b + \mathbf{s}, \end{cases} \quad (1)$$

où  $\mathbf{x}_b \in \mathbb{R}^n$  désigne une image de "fond", soulignant que, quel que soit le signal, l'image  $\mathbf{x}$  peut toujours être décomposée en un fond et un signal résiduel  $\mathbf{s} \in \mathbb{R}^n$ .

### 2.1 Observateur Idéal

Face à ces deux hypothèses (1), l'*observateur idéal* (IO) correspond au test statistique  $\Lambda_{io}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  qui minimise la probabilité moyenne d'erreur. Il est défini comme le rapport de vraisemblance entre les deux hypothèses [15] :

$$\Lambda_{io}(\mathbf{x}) = \frac{p(\mathbf{x}|H_1)}{p(\mathbf{x}|H_0)}. \quad (2)$$

Il est *optimal* car il établit une borne supérieure sur la performance de détection de tout observateur sur cette tâche.

### 2.2 Observateur Idéal Appris

L'*observateur idéal* étant analytiquement intractable, l'approche par apprentissage consiste à observer que toute transformation monotone du rapport de vraisemblance  $\Lambda_{io}(\mathbf{x})$  est également un test statistique optimal. Comme le note Zhou *et al.* [17], la probabilité postérieure  $\mathbb{P}(H_1|\mathbf{x})$  est une telle transformation :

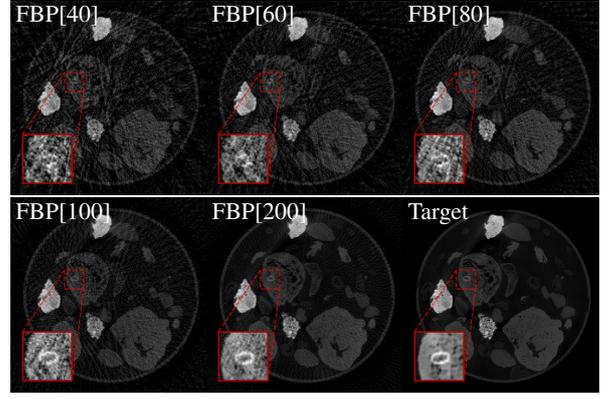


Figure 2 – **2DeteCT** : Disparition progressive du signal binaire synthétique en fonction du nombre de vues.

$$\mathbb{P}(H_1|\mathbf{x}) = \frac{(\mathbb{P}(H_1)/\mathbb{P}(H_0))\Lambda(\mathbf{x})}{1 + (\mathbb{P}(H_1)/\mathbb{P}(H_0))\Lambda(\mathbf{x})}. \quad (3)$$

Cela signifie qu'en utilisant (3), l'observateur idéal peut être appris en approximant la probabilité postérieure  $\mathbb{P}(H_1|\mathbf{x})$  sur une distribution donnée. Soit  $y \in \{0, 1\}$  la variable aléatoire associée à chaque hypothèse, et  $\mathbf{x}$  la variable aléatoire associée aux données, on peut montrer que minimiser la divergence de Kullback-Leibler entre la véritable probabilité postérieure et la postérieure approximée est équivalent à minimiser l'entropie croisée binaire (BCE) :

$$\begin{aligned} & \min_{\psi} \text{KL}(\mathbb{P}(y|\mathbf{x}) \parallel \mathbb{P}(y|\mathbf{x}, \psi)) \\ &= \min_{\psi} \mathbb{E}_{(\mathbf{x}, y)} [-y \log \Lambda_{\psi}(\mathbf{x}) - (1 - y) \log (1 - \Lambda_{\psi}(\mathbf{x}))], \end{aligned} \quad (4)$$

ici  $\Lambda_{\psi}(\mathbf{x})$  avec des paramètres  $\psi$  est la fonction approximant la probabilité postérieure de  $H_1$  sachant  $\mathbf{x}$ , *i.e.*  $\mathbb{P}(y = H_1|\mathbf{x}, \psi)$ . Naturellement,  $\Lambda_{\psi}$  peut être paramétré par un réseau de neurones profond, afin de pouvoir réutiliser des architectures efficaces [17] développées pour des tâches de classification.

## 3 Méthodologie

### 3.1 Problème Inverse

Dans ce travail, nous considérons le problème de tomographie par rayons X en vues éparses. Étant donné un opérateur de projection discret  $\mathbf{A} \in \mathbb{R}^{m \times n}$  et des mesures de projection  $\mathbf{b} \in \mathbb{R}^m$ , l'objectif est de trouver la véritable image  $\mathbf{x}^{\dagger} \in \mathbb{R}^n$  telle que

$$\mathbf{b} = \mathbf{A}\mathbf{x}^{\dagger} + \varepsilon,$$

avec  $\varepsilon \in \mathbb{R}^m$  le bruit de mesure. L'opérateur de rétro-projection filtrée (FBP)  $\mathbf{A}^+$  est utilisé pour générer des reconstructions bruitées  $\tilde{\mathbf{x}} = \mathbf{A}^+\mathbf{b}$ , ensuite post-traitées par un réseau de neurones profond  $D_{\phi}$  pour obtenir une estimation de la vérité terrain  $\hat{\mathbf{x}} = D_{\phi}(\tilde{\mathbf{x}})$ . Comme illustré dans les Figures 1 and 2, nous simulons des conditions d'imagerie de qualité différente en faisant varier le nombre de vues  $k > 0$  utilisées pour reconstruire les images. On suppose que les images cibles  $\mathbf{x}^{\dagger}$  sont échantillonnées à partir d'une distribution  $\pi_{\mathbf{x}^{\dagger}}$  et que les images d'entrée  $\tilde{\mathbf{x}}_k$  sont échantillonnées à partir d'une distribution  $\pi_{\tilde{\mathbf{x}}_k}$ , où l'indice  $k$  désigne le nombre de vues utilisées pour reconstruire l'image  $\tilde{\mathbf{x}}_k$ . nous supposons également que les images post-traitées  $\hat{\mathbf{x}}_k$  sont échantillonnées à partir d'une autre distribution  $\pi_{\hat{\mathbf{x}}_k}$ .

### 3.2 Détection de Signal Binaire

Pour tester cette méthode d'évaluation par observateur, nous utilisons le jeu de données **2DeteCT** [7]. Les mesures brutes sont agrégées afin de simuler un détecteur de 478 cellules, et des reconstructions de  $512^2$  pixels. Nous créons une tâche synthétique de détection de signal binaire. On fait l'hypothèse qu'une image  $\mathbf{x}$  contient un signal  $\mathbf{s}$ , qui est la réalisation d'une variable aléatoire  $s$ , et on le modélise comme une forme géométrique simple, *i.e.* une ellipse, avec une taille, une orientation et une position aléatoires dans l'image (Fig. 2).

### 3.3 Configurations d'Apprentissage

**Post-traitement** Pour le réseau de post-traitement, nous considérons une configuration simple, avec un réseau de type UNet [12] entraîné à prédire la vérité terrain  $\mathbf{x}^\dagger$  à partir de la reconstruction FBP en vues éparses  $\tilde{\mathbf{x}}$ . Nous considérons un opérateur général, qui prend également  $k$ , le nombre de vues, comme entrée, *i.e.*  $D_\phi(\cdot, k)$ , qui est entraîné à prédire l'image propre à partir de n'importe quelle configuration de vues éparses :

$$\min_{\phi} \mathbb{E}_k \mathbb{E}_{\mathbf{x}_k, \mathbf{x}} [\|D_\phi(\tilde{\mathbf{x}}_k, k) - \mathbf{x}^\dagger\|_2^2]. \quad (5)$$

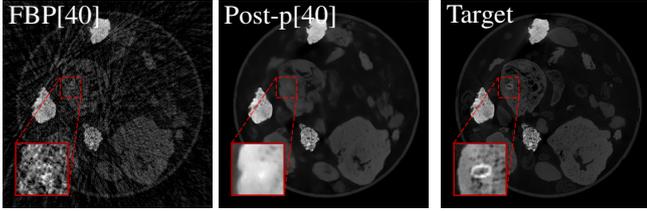


Figure 3 – **2DeteCT**: illustrations de reconstructions à 40 vues avec signal. Sur la reconstruction FBP, le signal est à peine visible, tandis que l'opération de post-traitement lisse complètement la région qui l'entoure.

**Apprentissage de l'observateur** Nous insistons sur le fait que pour approximer correctement l'observateur idéal (4), l'opérateur appris doit minimiser l'entropie croisée binaire pour la distribution d'images sur laquelle il est évalué. Cela signifie qu'utiliser un classifieur entraîné sur des images propres donne des résultats sous-optimaux lorsqu'il est évalué sur des échantillons FBP en vues éparses ou post-traités (voir Fig. 4–*gauche*). Par conséquent, pour maximiser la performance de l'*observateur* appris sur chaque cas, nous décidons d'entraîner un opérateur différent pour chaque configuration d'imagerie, c'est-à-dire que nous entraînons un classifieur spécialisé pour chaque nombre de vues  $k$  et pour chaque algorithme de reconstruction (voir Fig. 4–*droite*). Une approche moins coûteuse consisterait à *ajuster* un classifieur entraîné sur images propres pour chaque nouvelle configuration d'imagerie, mais nous ne l'explorons pas dans cet article.

**Optimisation relative à la tâche** Nous considérons l'objectif de minimisation dans (5) avec une légère modification, l'entropie croisée binaire est ajoutée comme terme de régularisation avec un paramètre  $\lambda > 0$

$$\min_{\phi} \mathbb{E}_k \mathbb{E}_{(\tilde{\mathbf{x}}_k, \mathbf{x}^\dagger, \mathbf{y})} \left[ \|D_\phi(\tilde{\mathbf{x}}_k, k) - \mathbf{x}^\dagger\|_2^2 - \lambda \log \mathbb{P}(\mathbf{y} | D_\phi(\tilde{\mathbf{x}}_k, k), \hat{\psi}) \right] \quad (6)$$

où  $\hat{\psi}$  sont les paramètres de l'*observateur* entraîné sur des images propres.

### 3.4 Configuration d'évaluation

Pour évaluer la performance d'une configuration d'imagerie spécifique, nous calculons l'erreur quadratique moyenne (MSE) ainsi que l'AUC[ROC] de son *observateur idéal* appris (voir Fig. 4).

## 4 Résultats

### 4.1 Approximer l'observateur idéal

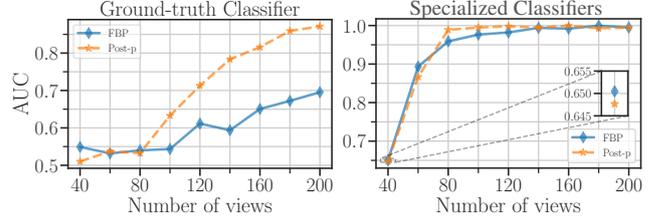


Figure 4 – **2DeteCT**: (*gauche*) évolution de l'AUC calculée par un classifieur entraîné sur des images propres, évalué sur des images bruitées  $\tilde{\mathbf{x}}_k$  et des images post-traitées  $\hat{\mathbf{x}}_k$ . (*droite*) évolution de l'AUC calculée par des classifieurs spécialisés.

Dans la section suivante, nous nous référerons interchangeablement à l'AUC comme la métrique *basée sur la tâche*. Relativement à notre problème initial Figure 1, on note que sur la Fig. 4, l'AUC classe correctement les reconstructions FBP avec 200 vues comme meilleures qu'une reconstruction post-traitée avec 40 vues.

Nous analysons maintenant si, pour un nombre donné de vues, l'opération de post-traitement améliore réellement les performances en classification. De manière remarquable, malgré l'amélioration de la MSE d'un ordre de grandeur (Fig. 1), l'AUC à 40 et 60 vues diminue après l'opération de post-traitement. Il apparaît alors que le signal soit si faible que l'opération de post-traitement supprime le signal de l'image, ou du moins diminue son intensité. Cela est confirmé par les résultats qualitatifs dans Fig. 3 où nous voyons que le post-traitement avec 40 vues lisse complètement l'image autour du signal. Pour un nombre plus élevé de vues, on remarque qu'avec un observateur correctement entraîné, le gain de performance apporté par le post-traitement de la reconstruction FBP reste limité Fig. 4. Nous faisons cette première observation et émettons l'hypothèse que maximiser le PSNR d'une reconstruction peut avoir un impact limité sur la réussite de détection en aval. Une approche plus spécifique consiste à optimiser le système d'imagerie, y compris l'algorithme de reconstruction, en fonction de la tâche en aval, ici la classification. C'est l'idée derrière l'optimisation *basée sur la tâche* (6).

### 4.2 Optimisation basée sur la tâche

Nous complétons notre analyse en comparant les performances de l'algorithme de reconstruction *basé sur la tâche* avec la méthode de post-traitement standard. Comme attendu, la régularisation de l'algorithme de reconstruction améliore les performances en classification (Fig. 5).

Nous fournissons également des résultats qualitatifs pour illustrer l'impact de l'optimisation *basée sur la tâche* sur la visibilité du signal pour de faibles nombres de vues.(Fig. 6).

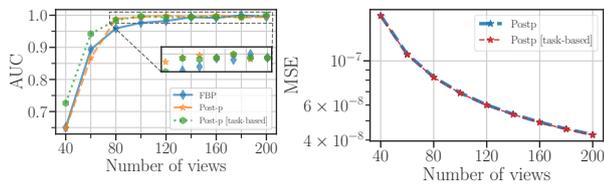


Figure 5 – **2DeteCT**: nous complétons les résultats Figure 4 en comparant les performances de l’algorithme de reconstruction basé sur la tâche (courbe verte). Nous montrons l’AUC de l’observateur appris, et la MSE des reconstructions.

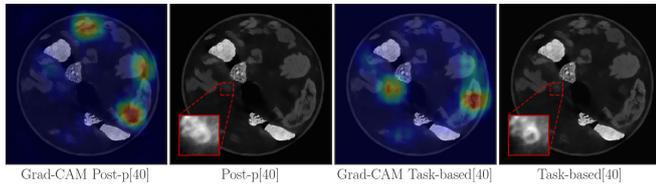


Figure 6 – **2DeteCT**: nous illustrons les améliorations apportées par l’optimisation basée sur la tâche de l’algorithme de reconstruction. Nous affichons les résultats à 40 vues et complétons l’illustration avec des cartes de chaleur Grad-CAM [13] pour interpréter le changement de performance en classification. Ici, l’observateur “remarque” la présence d’un signal - nous interprétons l’amélioration de l’AUC comme une diminution de l’entropie du classifieur.

## 5 Conclusion

Nous avons présenté une méthodologie d’évaluation par observateur pour quantifier la qualité d’une image relativement à une tâche en aval. Nous avons vu que la performance de la tâche reflète correctement la qualité globale des images produites. Nous avons montré comment le développement (et l’optimisation) d’un algorithme de reconstruction devait s’effectuer relativement à la tâche pour laquelle il est déployé. En ce sens, nous avons observé que les algorithmes basé-tâche donnaient les meilleurs résultats, validant effectivement l’approche. L’étape suivante serait d’évaluer l’approche par observateur sur des algorithmes plus avancés, *e.g.* méthodes déroulées, *plug-and-play* ou *diffusion*, afin d’observer si l’optimisation basée sur la tâche est bien déterminante pour la qualité de l’image ou si des algorithmes plus complexes peuvent compenser l’absence de régularisation par la tâche.

## References

- [1] Harrison H. BARRETT, Jie YAO, Jannick P. ROLLAND et Kyle J. MYERS : Model observers for assessment of image quality. *In National Academy of Sciences*, 1993.
- [2] Anna BREGER, Ander BIGURI, Malena Sabaté LANDMAN, Ian SELBY, Nicole AMBERG, Elisabeth BRUNNER, Janek GRÖHL, Sepideh HATAMIKIA, Clemens KARNER, Lipeng NING, Sören DITTMER, Michael ROBERTS, AIX-COVNET COLLABORATION et Carola-Bibiane SCHÖNLIEB : A study of why we need to reassess full reference image quality assessment with medical images. *Journal of Imaging Informatics in Medicine*, 2025.
- [3] Sebastian ECKEL, Peter HUTHWAITE, Michael LOWE, Andreas SCHUMM et Pierre GUÉRIN : Establishment and

validation of the Channelized Hotelling Model Observer for image assessment in industrial radiography. *NDT & E International*, 2018.

- [4] Yo Seob HAN, Jaejun YOO et Jong Chul YE : Deep Residual Learning for Compressed Sensing CT Reconstruction via Persistent Homology Analysis. 2016.
- [5] Zohaib Amjad KHAN, Giuseppe VALENZISE, Aladine CHETOUANI et Frédéric DUFAUX : Towards an image utility assessment framework for machine perception. *In European Signal Processing Conference*, 2022.
- [6] Maximilian B. KISS, Sophia B. COBAN, K. Joost BATENBURG, Tristan van LEEUWEN et Felix LUCKA : 2DeteCT - A large 2D expandable, trainable, experimental Computed Tomography dataset for machine learning. *Scientific Data*, 2023.
- [7] Matthew A. KUPINSKI, Darrin C. EDWARDS, Maryellen L. GIGER et Charles E. METZ : Ideal observer approximation using Bayesian classification neural networks. *IEEE Transactions on Medical Imaging*, 2001.
- [8] Matthew A. KUPINSKI, John W. HOPPIN, Eric CLARKSON et Harrison H. BARRETT : Ideal-observer computation in medical imaging with use of Markov-chain Monte Carlo techniques. *Journal of the Optical Society of America*, 2003.
- [9] Kaiyan LI, Weimin ZHOU, Hua LI et Mark A. ANASTASIO : Assessing the Impact of Deep Neural Network-Based Image Denoising on Binary Signal Detection Tasks. *IEEE Transactions on Medical Imaging*, 2021.
- [10] Robert N. MCDONOUGH, Anthony D. WHALEN et Anthony D. WHALEN : *Detection of Signals in Noise*. Academic Press, 2nd ed édition, 1995.
- [11] Olaf RONNEBERGER, Philipp FISCHER et Thomas BROX : U-Net: Convolutional Networks for Biomedical Image Segmentation. *In Medical Image Computing and Computer-Assisted Intervention*, 2015.
- [12] Ramprasaath R. SELVARAJU, Michael COGSWELL, Abhishek DAS, Ramakrishna VEDANTAM, Devi PARIKH et Dhruv BATRA : Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization. 2020.
- [13] Zhou WANG, Alan Conrad BOVIK, Hamid Rahim SHEIKH et Eero P. SIMONCELLI : Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 2004.
- [14] Anthony WHALEN : Statistical theory of signal detection and parameter estimation. *IEEE Communications Magazine*, 1984.
- [15] Xiaohui ZHANG, Varun A. KELKAR, Jason GRANSTEDT, Hua LI et Mark A. ANASTASIO : Impact of deep learning-based image super-resolution on binary signal detection. *Journal of Medical Imaging*, 2021.
- [16] Weimin ZHOU, Hua LI et Mark A. ANASTASIO : Approximating the Ideal Observer and Hotelling Observer for binary signal detection tasks by use of supervised learning methods. *IEEE Transactions on Medical Imaging*, 2019-10.