

# La phylogénie des images dans les réseaux sociaux

---



## Étude bibliographique

Master *Sciences et Technologies*,  
Mention *Informatique*,  
Parcours IMAGINA

**Auteur**

Noé Le Philippe

**Superviseur**

William Puech

**Lieu de stage**

Équipe ICAR - LIRMM UM5506 - CNRS, Université de Montpellier

## Résumé

Ce stage de master.

---

## Abstract

This master thesis.

---

# Table des matières

<b>Table des matières</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Near-duplicate images (NDI) . . . . .	2
1.2 Arbre phylogénétique (Image Phylogeny Tree - IPT) . . . . .	2
1.3 Pourquoi se restreindre à la compression ? . . . . .	3
<b>2 État de l'art</b>	<b>4</b>
2.1 Étude de l'arbre phylogénétique . . . . .	4
2.2 Analyse des recompression JPEG . . . . .	6
2.3 Distances entre distributions . . . . .	8
<b>3 Notre approche</b>	<b>9</b>
3.1 Principe . . . . .	9
3.2 La construction de l'arbre . . . . .	9
<b>4 Conclusion</b>	<b>12</b>
<b>Bibliographie</b>	<b>13</b>

## Introduction

La phylogénie, en sciences naturelles, est définie[12] comme l'étude des relations de parenté entre êtres vivants. Et c'est exactement de cela qu'il s'agit dans le cas des images, l'étude des relations de parenté entre images.

À l'ère des réseaux sociaux, il n'a jamais été aussi simple de partager des idées et du contenu. À chaque partage cependant, l'information peut être amenée à être modifiée. Les images, puisque c'est là notre sujet d'étude, peuvent avoir subi un certain nombre de transformations et modifications avant de nous parvenir. C'est dans ce contexte que nous allons intervenir, et tenter de reconstituer la phylogénie de l'image. Il peut être difficile de différencier cette image de l'originale, et de savoir laquelle est l'originale, mais c'est pourtant crucial dans un monde où l'information peut être falsifiée par tout le monde et extrêmement facilement. Les applications sont variées, et ne se cantonnent pas à la détection et la discrimination d'images altérées, on peut également se servir de la phylogénie de l'image pour optimiser de l'espace de stockage en ne gardant que l'originale ou encore suivre la diffusion et l'évolution des idées sur les réseaux sociaux.



FIGURE 1.1 – Exemples de near-duplicates

## 1.1 Near-duplicate images (NDI)

Nous travaillons sur un ensemble d'images, toutes similaires, et au milieu de cette jungle d'images, nous devons décider quelle image est le parent de quelle autre, ou autrement dit, quelles images sont des **near-duplicates**. Joly et al. [7] définit la notion de near-duplicate comme suit :  $I_1 = T(I)$ ,  $T \in \mathcal{T}$  où  $I$  est l'image parent,  $I_1$  est l'image enfant et  $\mathcal{T}$  est un ensemble de transformations autorisées,  $I_1$  et  $I$  sont alors des NDI. Dans le cas général,  $\mathcal{T} = \{\text{resampling, cropping, affine warping, color changing, lossy compression}\}$ , Fig. 1.1 montre un exemple de near-duplicates, dans le cadre de ce stage cependant,  $\mathcal{T} = \{\text{lossy compression}\}$ . Notons l'utilisation du terme *transformations autorisées*. Ce terme est double, d'une part, il place une limite arbitraire dans la force de la transformation, par exemple une image croquée à plus de 10% pourra ne pas pas être considérée comme un near-duplicate, et d'autre part, il permet de restreindre l'espace des transformations possibles. Ainsi, dans le cadre du stage, seules les images de la troisième case de Fig. 1.1 seront des NDI. Ces transformations peuvent évidemment se composer, et une image enfant peut être le résultat de plusieurs transformations.

**Note.** Nous utiliserons relation parent-enfant et relation de parenté de manière interchangeable dans le reste de ce rapport, mais c'est bien d'une relation parent-enfant qu'il s'agit.

## 1.2 Arbre phylogénétique (Image Phylogeny Tree - IPT)

C'est l'arbre représentant les relations de parenté entre les différentes images. Il sera extrait d'un ensemble de NDI, c'est l'objectif final de l'application. Un exemple est disponible Fig. 1.2. Le passage d'une génération à l'autre, autrement dit d'un noeud à son fils se fait à travers la transformation  $I_1 = T(I)$ , ainsi, une image et son parent sont des NDI, alors qu'une images et ses soeurs ne le sont pas (Fig 1.3).

La reconstruction de l'arbre se concentre autour de deux problèmes principaux. Le premier est l'identification de la racine, et le second et l'estimation du reste de l'arbre. Il est en effet critique d'identifier correctement la racine. Prenons par exemple un des cas d'utilisation de l'IPT, la détection d'altération d'images. L'idée est que pour un ensemble d'images, plus on est proche de la racine, moins l'image a subi de transformations, et donc moins elle est altérée, avec dans le meilleur des cas, la racine en image originale. On voit bien que si on identifie mal la racine, on déduira, à tort, qu'une image n'a pas été altérée. Il n'est pas toujours garanti que la totalité des

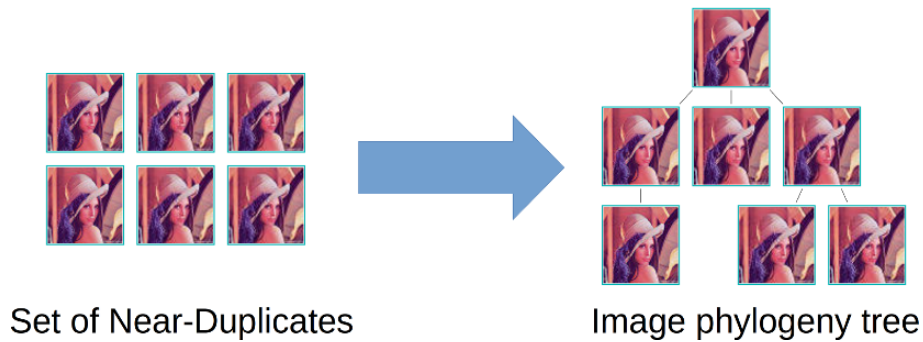


FIGURE 1.2 – Passage d'un ensemble de NDI à un arbre phylogénétique



FIGURE 1.3 – Passage d'une image parent à l'image enfant

images de l'arbre original soit présente, de plus, certaines transformations peuvent être mineures, et difficile à détecter, c'est donc bien une estimation de l'arbre original qui sera faite.

### 1.3 Pourquoi se restreindre à la compression ?

Comme mentionné précédemment, nous ne considérons que le cas de la compression avec perte en laissant de côté les autres transformations. Nous avons décidé de ne nous concentrer que sur une seule transformation pour avoir la possibilité de la traiter en détail et en profondeur dans le cadre du stage et ne pas devoir survoler sans approfondir toutes les transformations. Le choix de la compression est assez naturel, cela n'a, à notre connaissance, pas encore été traité dans le cadre de la phylogénie, et c'est un domaine largement étudié en forensic.

## État de l'art

### 2.1 Étude de l'arbre phylogénétique

#### La Visual Migration Map (VMM)

C'est à notre connaissance, le premier article concernant vraiment notre sujet. Kennedy et al. [9] proposent une approche permettant d'automatiquement détecter la manière dont une image a été éditée ou manipulée, et d'en extraire des relations parent-enfant entre les images. Il vont construire à partir de l'estimation de ces transformations une Visual Migration Map (VMM) (voir Fig. 2.1) qui est en fait notre arbre de phylogénie.

Ils partent du principe que les transformations sont directionnelles, c'est à dire que l'on ne peut passer que d'une image moins transformée à une image plus transformée. Ainsi, ils vont tenter d'estimer la direction de chaque transformation entre deux images  $I_1$  et  $I_2$  (sachant que  $\mathcal{T} = \{\textit{scaling}, \textit{cropping}, \textit{grayscale}, \textit{overlay}, \textit{insertion}\}$ ). Trois scénarios sont alors possibles : si toutes les transformations sont dans le même sens, l'image fille est alors celle vers qui pointent les transformations, si les transformations sont dans des sens contraires, les images sont sûrement des soeurs, elles n'ont en tous cas pas de relation parent-enfant, et enfin si aucune transformation n'a été détectée, c'est que soit les images sont identiques, soit qu'elles ne sont pas des near-duplicates. On peut en voir un exemple Fig. 2.2.

Un graphe va ensuite être construit à partir des couples d'images pour lesquels une relation parent-enfant a été détectée. À noter qu'une relation parent-enfant ne veut pas forcément dire que c'est le parent direct mais plutôt un ancêtre. Ainsi, un noeud du graphe (une image) peut avoir plusieurs noeuds parents, pour finalement obtenir l'arbre désiré, seuls les chemins les plus longs sont conservés, comme on peut le voir Fig. 2.3.

#### Image Phylogeny Tree (IPT)

Tout comme l'approche présentée précédemment, Dias et al. [3][4] proposent une approche basée sur le contenu de l'image. Cela consiste à mapper une image sur le domaine d'une autre, pour pouvoir les comparer, et ensuite estimer le coût de cette opération, cela fait comme hypothèse que si deux images sont dépendantes, alors il est possible d'obtenir l'une en appliquant une transformation à l'autre.

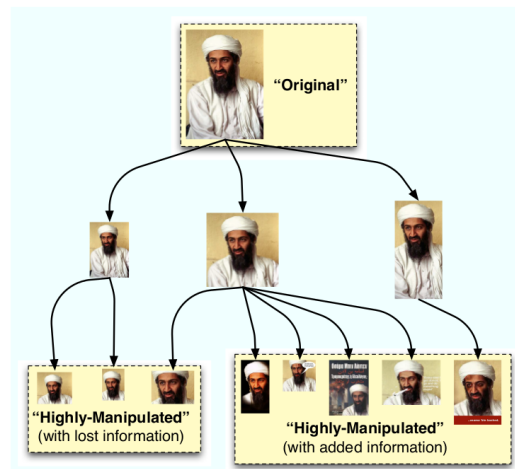


FIGURE 2.1 – Exemple de VMM, issu de [9]

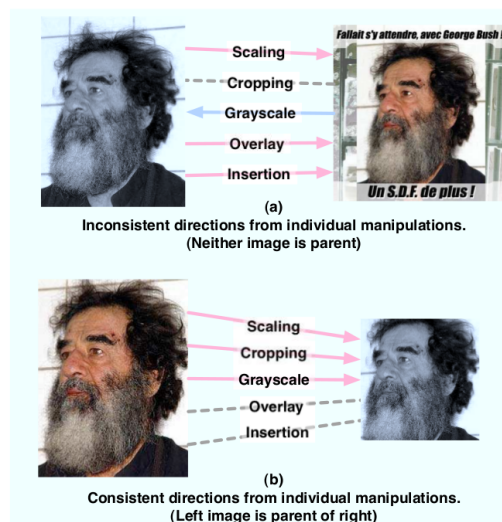


FIGURE 2.2 – Exemple de direction des transformations, issu de [9]

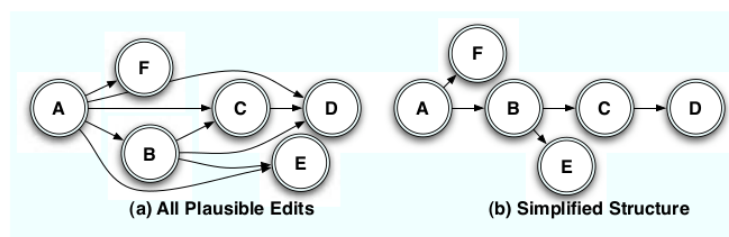


FIGURE 2.3 – Exemple de simplification de graphe, issu de [9]



Les images sont comparées à l'aide d'une *dissimilarity function*  $d(.,.)$  qui renvoie des petites valeurs lorsque les images sont proches (elles ont subi des transformations similaires). L'équation 2.1 détaille cette fonction.  $A$  et  $B$  sont les deux images qui vont être comparées, et  $T_{\vec{\beta}} \in \mathcal{T}$  est la transformation en train d'être estimée. La transformation la plus faible est gardée comme résultat, c'est la transformation la plus probable qu'a pu subir l'image. À noter que  $d(.,.)$  est asymétrique, c'est à dire que  $d(A, B) \neq d(B, A)$ , ce qui est parfaitement logique, en plus d'être nécessaire, sachant que les transformation sont directionnelles, comme expliqué précédemment.

$$d(A, B) = \min_{T_{\vec{\beta}}} \left| B - T_{\vec{\beta}}(A) \right|_{\text{comparison method}} \quad (2.1)$$

On voit donc bien que la problématique principale de cette méthode est de trouver une bonne méthode de comparaison. On pourra toutefois regretter que les auteurs soient restés assez vagues sur leurs méthodes de comparaison.

$d(.,.)$  est donc appliquée à tous les couples d'images de l'ensemble, pour créer une *dissimilarity matrix*, une matrice de taille  $n \times n$  qui sera ensuite donnée comme entrée à leur algorithme de reconstruction d'arbre, Oriented Kruskal. Nous n'allons pas entrer dans les détails de cet algorithme, il est expliqué de manière exhaustive dans [4], et ayant une approche différente, comme détaillé Chap. 3, la reconstruction sera également différente.

En plus de proposer une approche pour reconstituer l'arbre de phylogénie, ils proposent une approche pour comparer deux arbres, et donc évaluer notre arbre reconstruit si on le compare avec la vérité terrain. Cela consiste en quatre métriques :

$$\mathbf{Root} \quad R(IPT_1, IPT_2) = \begin{cases} 1 & \text{if } \text{Root}(IPT_1) = \text{Root}(IPT_2) \\ 0 & \text{Otherwise} \end{cases}$$

$$\mathbf{Edges} \quad E(IPT_1, IPT_2) = \frac{|E_1 \cap E_2|}{n-1}$$

$$\mathbf{Leaves} \quad L(IPT_1, IPT_2) = \frac{|L_1 \cap L_2|}{|L_1 \cup L_2|}$$

$$\mathbf{Ancestry} \quad A(IPT_1, IPT_2) = \frac{|A_1 \cap A_2|}{|A_1 \cup A_2|}$$

**Root** est triviale, et renvoie si les racines sont identiques. **Edges** mesure le ratio de noeuds ayant le bon parent, **Leaves** est le ratio de feuilles correctes, et enfin **Ancestry** est le ratio d'ancêtres corrects jusqu'à la racine.

Ces métriques serviront à évaluer nos résultats dans la suite du stage.

## 2.2 Analyse des recompression JPEG

### Qu'est ce que JPEG ?

Une rapide introduction sur ce qu'est le JPEG et son fonctionnement s'impose, nous ne parlerons cependant que du mode de compression avec perte, en laissant de côté son mode de compression sans perte, peu intéressant en plus de n'être virtuellement jamais utilisé. Pour de plus amples détails, le lecteur se dirigera vers [13].

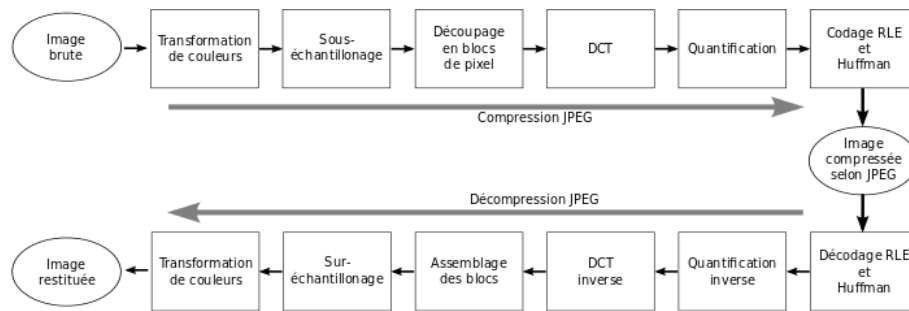


FIGURE 2.4 – Étapes de compression et décompression, issu de [8]

Fig. 2.4 liste toutes les étapes permettant de passer d'une image non compressée à une image compressée pour finir par la décompresser, mais les étapes sont sommairement :

- L'image est convertie dans l'espace YUV
- Les canaux U et V sont sous-échantillonnés
- Chaque canal est découpé en blocs de  $8 \times 8$
- Une DCT est appliquée à chacun de ces blocs
- Chaque coefficient est quantifié selon la table de quantification (voir Fig. 2.5) correspondant au facteur de qualité  $Q \in \{1, 2, 3, \dots, 100\}$
- Le tout est ensuite compressé à l'aide d'un codage entropique

C'est surtout l'étape de quantification qui va réduire la quantité d'information, et permettre de réduire la taille du fichier. Cette quantification, si elle est trop agressive va faire apparaître des artefacts de bloc, des blocs  $8 \times 8$  visibles (voir Fig. 2.6). C'est ce qui caractérise le JPEG, et qui permet de détecter un certain nombre de choses, comme l'altération d'image [2], les doubles compressions [1] ou encore dans notre cas les relations de parenté.

$$Q = \begin{bmatrix} 16 & 11 & 10 & 16 & 24 & 40 & 51 & 61 \\ 12 & 12 & 14 & 19 & 26 & 58 & 60 & 55 \\ 14 & 13 & 16 & 24 & 40 & 57 & 69 & 56 \\ 14 & 17 & 22 & 29 & 51 & 87 & 80 & 62 \\ 18 & 22 & 37 & 56 & 68 & 109 & 103 & 77 \\ 24 & 35 & 55 & 64 & 81 & 104 & 113 & 92 \\ 49 & 64 & 78 & 87 & 103 & 121 & 120 & 101 \\ 72 & 92 & 95 & 98 & 112 & 100 & 103 & 99 \end{bmatrix}$$

FIGURE 2.5 – Exemple de table de quantification, issu de [8]

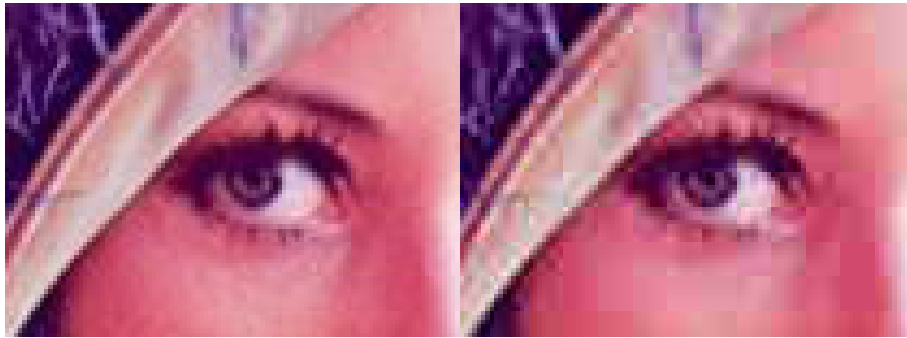


FIGURE 2.6 – Artefacts de blocs, à gauche avec  $Q=90$ , à droite avec  $Q=20$

### Détection des recompressions JPEG

[5]

### Détection des double compressions JPEG avec la même matrice de quantification

[6]

### Convergence des blocs lors de compressions successives

[10]

## 2.3 Distances entre distributions

[11]

## Notre approche

### 3.1 Principe

**Théorème.** *Pour tout couple d'image  $(A, B)$ , s'il n'est pas possible de prouver que  $A$  n'est pas le parent de  $B$ , alors il y a une relation parent-enfant entre  $A$  et  $B$ ,  $A \rightarrow B$ .*

*Démonstration.* Soit  $f(A, B)$  une fonction qui pour tout couple d'image  $(A, B)$  détecte à chaque qu'il est présent un marqueur prouvant qu'il n'y a pas de relation de parenté entre  $A$  et  $B$ . Si  $f(A, B)$  ne détecte rien alors c'est que  $A$  est le parent de  $B$ .  $\square$

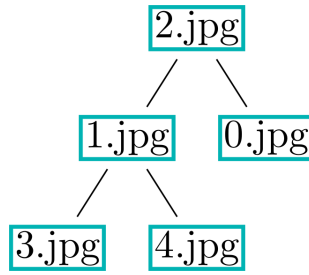
### 3.2 La construction de l'arbre

On va ainsi pour tout couple d'image de notre ensemble tenter de nier qu'il existe une relation de parenté, cela va permettre d'extraire une matrice binaire, de la taille de l'ensemble d'image, où une case à vrai indiquera une relation de parenté. À noter une fois de plus que ce n'est pas forcément un parent direct mais plutôt un ancêtre.

On peut noter plusieurs choses intéressantes de l'exemple Fig. 3.1. L'image  $I_2$  a toute sa colonne marquée à vrai, c'est à dire que c'est un parent commun à toutes les autres images, et sa ligne n'a que des 0, elle n'a donc aucun parent. On se sert de ce principe pour la reconstruction de l'arbre à partir de la matrice : une image qui n'a pas de parent est la racine. On peut également remarquer les colonnes où toutes les cases sont à 0, cela veut dire que ces images ne sont le parent de personne, et donc sont des feuilles.

### L'algorithme de reconstruction

De ces quelques principes nous avons extrait un algorithme.



(a) Arbre de phylogénie

-	$I_0$	$I_1$	$I_2$	$I_3$	$I_4$
$I_0$	-	0	0	0	0
$I_1$	0	-	0	1	1
$I_2$	1	1	-	1	1
$I_3$	0	0	0	-	0
$I_4$	0	0	0	0	-

(b) Matrice de parenté

FIGURE 3.1 – Une arbre de phylogénie et sa matrice de parenté

**Data:**  $M$  a  $n \times n$  parentage matrix

**Result:** the root of the tree

```

1  $nextRoot \leftarrow$  row with min sum of elements;
2  $treeRoot \leftarrow nextRoot$ ;
3 forall the rows  $row$  of  $M$  do
4    $root \leftarrow nextRoot$ ;
5   mark  $root$  as done;
6   for  $i \leftarrow 0$  to  $n$  do
7      $row[i] \leftarrow 0$ ;
8     if sum of elements of  $row == 0$  then
9       | add  $i$  as child of  $root$ ;
10    end
11    if  $row$  has the smallest sum of elements and is not marked as done then
12      |  $nextRoot \leftarrow i$ ;
13    end
14  end
15 end
16 return  $treeRoot$ 

```

Il prend en entrée la matrice de parenté détaillée précédemment et retourne la racine de l'arbre. La philosophie de l'algorithme est qu'une image n'ayant aucun parent est la racine, et ce de manière récursive grâce à la structure d'arbre.

À chaque tour de boucle, une image est sélectionnée comme la racine (lignes 1 et 11) et nommée  $root$ , elle est retirée (ligne 7) des ancêtres des autres images si c'est un ancêtre. Si ces autres images n'ont plus d'ancêtre (ligne 8), c'est que  $root$  était le parent direct de l'image en train d'être traitée, cette image est donc ajoutée comme enfant de  $root$  (ligne 9). La ligne 5 permet de ne traiter qu'une fois chaque image comme racine potentielle.

Cet algorithme a une complexité de  $O(n^2)$ . Il y a deux boucles imbriquées, et si les sommes sont calculées une seule fois au début et mises à jour à chaque fois qu'un parent est enlevé, il n'y a pas de boucle supplémentaire augmentant la complexité.

Le but est donc de trouver la fonction permettant de détecter le marqueur prouvant qu'il n'y a pas de parenté.

Il nous a semblé important, même dans le cadre de l'étude bibliographique, d'élaborer cet algorithme. Il est en effet assez simple mais néanmoins indispensable à notre méthode, et aurait pu orienter différemment nos recherches et notre méthode s'il n'avait pas été imaginé. Le reste du stage sera consacré à élaborer la fonction énoncée précédemment, et c'est là que la majorité du travail se concentrera.

Cette approche a pour avantage de réduire l'évaluation d'un arbre de phylogénie à partir d'un ensemble d'images à l'évaluation binaire de parenté entre deux images.

---

## Conclusion

---

## Bibliographie

- [1] Tiziano BIANCHI et Alessandro PIVA. “Detection of nonaligned double JPEG compression based on integer periodicity maps”. In : *Information Forensics and Security, IEEE Transactions on* 7.2 (2012), p. 842–848.
- [2] Tiziano BIANCHI et Alessandro PIVA. “Image forgery localization via block-grained analysis of JPEG artifacts”. In : *Information Forensics and Security, IEEE Transactions on* 7.3 (2012), p. 1003–1017.
- [3] Zanoni DIAS, Anderson ROCHA et Siome GOLDENSTEIN. “First steps toward image phylogeny”. In : *Information Forensics and Security (WIFS), 2010 IEEE International Workshop on*. IEEE. 2010, p. 1–6.
- [4] Zanoni DIAS, Anderson ROCHA et Siome GOLDENSTEIN. “Image phylogeny by minimal spanning trees”. In : *Information Forensics and Security, IEEE Transactions on* 7.2 (2012), p. 774–788.
- [5] Xiaoying FENG et Gwenaél DOËRR. “JPEG recompression detection”. In : *IS&T/SPIE Electronic Imaging*. International Society for Optics et Photonics. 2010, 75410J–75410J.
- [6] Fangjun HUANG, Jiwu HUANG et Yun Qing SHI. “Detecting double JPEG compression with the same quantization matrix”. In : *Information Forensics and Security, IEEE Transactions on* 5.4 (2010), p. 848–856.
- [7] Alexis JOLY, Olivier BUISSON et Carl FRÉLICOT. “Content-based copy retrieval using distortion-based probabilistic similarity search”. In : *Multimedia, IEEE Transactions on* 9.2 (2007), p. 293–306.
- [8] *JPEG*. URL : <https://fr.wikipedia.org/wiki/JPEG>.
- [9] Lyndon KENNEDY et Shih-Fu CHANG. “Internet image archaeology : automatically tracing the manipulation history of photographs on the web”. In : *Proceedings of the 16th ACM international conference on Multimedia*. ACM. 2008, p. 349–358.
- [10] ShiYue LAI et Rainer BOHME. “Block convergence in repeated transform coding : JPEG-100 forensics, carbon dating, and tamper detection”. In : *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE. 2013, p. 3028–3032.
- [11] Marina A OIKAWA et al. “Distances in multimedia phylogeny”. In : *International Transactions in Operational Research* (2015).



- [12] *Phylogénie*. URL : <https://fr.wikipedia.org/wiki/Phylogen%C3%A8se>.
- [13] Gregory K WALLACE. “The JPEG still picture compression standard”. In : *Consumer Electronics, IEEE Transactions on* 38.1 (1992), p. xviii–xxxiv.