

***Inpainting* vidéo pour la restauration de films par reconstructions alternées de la structure et de la texture**

Arthur RENAUDEAU¹ François LAUZE² Fabien PIERRE³ Jean-François AUJOL⁴ Jean-Denis DUROU¹

¹ IRIT, UMR CNRS 5505, Université de Toulouse

² DIKU, Université de Copenhague, Danemark

³ LORIA, UMR CNRS 7503, Université de Lorraine, INRIA projet Magrit

⁴ IMB, UMR CNRS 5251, Université de Bordeaux, Bordeaux INP

arthur.renaudeau@irit.fr

Résumé

Nous proposons un nouveau modèle d'inpainting vidéo pour la restauration de films, qui combine la reconstruction de la structure par une méthode de diffusion et la reconstruction de la texture par une méthode de recopie de patches. Les énergies proposées pour chacune de ces deux méthodes sont minimisées alternativement, afin de préserver la structure globale de l'image tout en affinant sa texture. Alors que la reconstruction de la structure est effectuée conjointement à l'estimation du mouvement par flux optique via plusieurs approches proximales, la reconstruction de la texture est traitée par une approche variationnelle non locale (NL-means). Les résultats sur différentes séquences d'images de la base de données Middlebury et de la Cinémathèque de Toulouse montrent une amélioration dans la qualité des reconstructions.

Mots Clef

Inpainting vidéo, flux optique, méthodes variationnelles.

Abstract

We propose a new video inpainting model for movies restoration application. This model combines structural reconstruction with a diffusion-based method and textural reconstruction with a patch-based method. Both proposed energies (one for each method) are alternatively minimized in order to preserve the overall structure adding textural refinement. While the structural reconstruction is obtained jointly with optical flow computation thanks to several proximal approaches, the textural reconstruction is processed by a variational non-local approach (NL-means). Results on different frames from the Middlebury database and from the Cinémathèque de Toulouse show quality improvement in the reconstructions.

Keywords

Video inpainting, optical flow, variational methods.

1 Introduction

L'*inpainting* vidéo est un problème clé de l'industrie cinématographique, qui peut aider à automatiser la restauration des films ayant subi des dégradations importantes (voir FIGURE 1), ou à mettre en œuvre des effets spéciaux nécessitant le retrait de certains éléments (« réalité diminuée »). L'*inpainting* vidéo, comme tout traitement vidéo, est de plus en plus utilisé grâce à l'augmentation des performances des processeurs et des GPU qui permet d'effectuer des calculs à grande échelle. Les techniques d'*inpainting* vidéo utilisent soit des méthodes de diffusion avec estimation du mouvement, soit des méthodes reposant sur la manipulation de patches 3D afin de prendre en compte la redondance temporelle entre images consécutives, mais sans estimation explicite du mouvement dans ce cas. Or, le mouvement pouvant donner beaucoup d'informations pour récupérer des données, il s'agit d'un indice précieux pour effectuer l'*inpainting*. Cependant, pour estimer le mouvement, il est nécessaire de disposer de données complètes. C'est pourquoi cette estimation doit être traitée en même temps que l'*inpainting* lui-même, ce qui représente notre principal objectif.



FIGURE 1 – Exemple d'une séquences de trois images numérisées issues d'un vieux film de la Cinémathèque de Toulouse, dont l'image centrale comporte un défaut.

Dans cet article, nous visons à restaurer les défauts de type « taches » sur des films déjà numérisés. Chacun de ces défauts apparaît dans une seule image, et non sur les images situées juste avant ou juste après dans la séquence. Pour les éliminer, notre approche consiste à combiner un modèle

d’*inpainting* vidéo par diffusion calculant conjointement le flux optique, avec un modèle utilisant des patches 2D et des cartes de correspondances vers les images voisines temporellement. Grâce à cette approche, notre modèle n’a donc besoin que des deux images adjacentes pour reconstruire la zone endommagée. Si chaque modèle pris séparément présente des inconvénients en termes de qualité de reconstruction, leur combinaison permet d’améliorer les résultats.

Après avoir décrit l’état de l’art dans la section 2, nous présentons notre modèle combiné dans la section 3. Notre stratégie numérique pour résoudre ce problème variationnel, présentée dans la section 4, consiste à optimiser alternativement les modèles utilisant la diffusion et les patches. Les résultats de notre méthode d’*inpainting* vidéo sur différentes séquences, présentés dans la section 5, confirment l’intérêt d’une telle combinaison de modèles.

2 État de l’art

L’*inpainting* est le nom donné à la technique qui consiste à remplir les parties endommagées ou manquantes d’une image. Le terme « *inpainting* » n’a été utilisé qu’à partir de 2000 dans [4], par analogie avec le processus de restauration utilisé dans le domaine de l’art, et après celui de « désoccultation » dans [21] en 1998. Les premières applications de l’*inpainting* proviennent à l’origine de modèles de diffusion pour le débruitage, qui remontent au début des années 1990. Ce champ de recherche a été très actif ces dernières années, stimulé par de nombreuses applications : suppression de rayures ou de texte superposé à une image, restauration d’une image altérée suite à une transmission, élimination d’objets dans le contexte de la réalité diminuée.

Le remplissage de la zone à restaurer est un problème inverse mal posé, car il n’existe pas de solution unique bien définie. Il est donc nécessaire d’introduire une connaissance a priori dans le modèle. Toutes les méthodes existantes sont guidées par l’hypothèse que les pixels situés dans les parties connues et manquantes de l’image partagent les mêmes propriétés statistiques ou structures géométriques. Cette hypothèse se retrouve dans différentes hypothèses a priori, locales ou globales, afin d’obtenir une image restaurée qui soit visuellement plausible.

L’*inpainting* par diffusion vise à propager l’information contenue dans les pixels, depuis le bord de la zone endommagée vers l’intérieur de celle-ci. La variation totale pour l’*inpainting* a été introduite dans [11] pour bloquer la diffusion sur les contours des objets et induire des restaurations constantes par morceaux. L’extension de l’*inpainting* par diffusion à la vidéo a commencé dans [12] et [17, 18], avec une estimation conjointe du mouvement pour guider le remplissage de la zone endommagée. À partir du modèle de flux optique bien connu de [16] avec une régularisation lisse L^2 , [2] a opté pour des normes L^1 afin de préserver les discontinuités du mouvement, en

effectuant la résolution avec les algorithmes proximaux de [23]. Très récemment, [7] a choisi un modèle TV- L^1 pour résoudre l’estimation du mouvement et la reconstruction d’images, en utilisant également des algorithmes proximaux.

Cependant, ces modèles par diffusion sont limités car ils ne peuvent pas traiter les textures. C’est pourquoi des modèles par recopie complète ou partielle de patches (détaillés dans [8]) ont été développés pour restaurer les détails, d’abord pour la synthèse de texture dans [15], puis pour l’*inpainting* local (la recherche de patches s’effectue dans un voisinage du défaut) dans [13] avec priorité au remplissage selon l’intensité du gradient spatial sur les bords de la zone à restaurer. Enfin, les méthodes récentes recourent à un mélange de patches suivant une recherche spatiale non locale comme dans [1]. Cette approche a également été étendue aux vidéos dans [22] et [19] avec des patches 3D, ce qui permet une recherche de similarité temporelle dans les comparaisons entre patches. Si ces modèles permettent de mieux récupérer les textures, ils sont cependant très dépendants de l’initialisation de la zone à remplir, afin de ne pas rester bloqué dans un minimum local suite aux différentes étapes de minimisation et entraîner une mauvaise reconstruction des structures régulières. De plus, la prise en compte de la taille du patch est également un critère important pour récupérer des textures ayant des propriétés statistiques différentes.

Dans le contexte de l’*inpainting* d’images, l’idée de mélanger les deux approches a déjà été étudiée dans [14], où la diffusion et le remplissage de la texture sont traités séquentiellement, dans [6] où l’image est décomposée en « cartoon » et « texture », qui sont restaurées séparément, et dans [9] où la reconstruction de la texture est guidée par les lignes de niveaux. Dans le contexte du débruitage vidéo, [5] utilise des patches combinés au calcul du flux optique de [23]. Notre but dans cet article est également de tirer le meilleur parti des deux approches, en combinant la diffusion pour récupérer la structure et les patches pour récupérer la texture, et en utilisant également les approches par diffusion et recopie de patches pour estimer le mouvement.

3 Énoncé du problème

Définissons une séquence de trois images couleur successives $\{u_b, u, u_f\}$ comme des fonctions $\Omega \subset \mathbb{R}^2 \rightarrow \mathbb{R}^3$ à variation bornée, où u_b est l’image arrière, u est l’image courante (ou centrale) contenant le défaut à corriger et u_f l’image avant. Cette zone de défaut est définie par $O \subset \Omega$. Le problème est alors le suivant :

$$u^* = \operatorname{argmin}_u \{E_S(v_b, u, v_f) + E_T(\Gamma_b, u, \Gamma_f)\} \quad (1)$$

où E_S l’énergie de reconstruction de la structure, minimisée canal par canal (ou avec la luminance pour l’estimation

du mouvement), et E_T celle de reconstruction de la texture, minimisée directement sur l'image couleur. Les différentes variables apparaissant dans (1) sont explicitées dans les prochaines sous-sections.

3.1 Énergie de reconstruction de la structure

Le modèle choisi pour E_S se fonde sur les travaux de [18] et [7], en ajoutant un traitement symétrique des flux optiques avant et arrière (rappelons que u comporte ici un seul canal) :

$$E_S(v_b, u, v_f) = \mu \int_{\Omega} |\nabla u(x)| dx + \lambda \int_{\Omega} |Jv_b(x)| dx + \lambda \int_{\Omega} |Jv_f(x)| dx + \int_{\Omega} |u_b(x + v_b(x)) - u(x)| dx + \int_{\Omega} |u_f(x + v_f(x)) - u(x)| dx \quad (2)$$

sous la contrainte $u = u^0$ sur $O^c = \Omega \setminus O$, afin de préserver la partie saine de l'image. Les termes faisant intervenir le champ de déplacements $v_b : \Omega \rightarrow \mathbb{R}^2$ (respectivement v_f) traduisent la contrainte de flux optique régularisée en norme L^1 , proposée par [23], entre l'image courante u et l'image arrière u_b (respectivement l'image avant u_f), mais en utilisant des matrices jacobiniennes, comme dans [10], pour conserver un caractère isotrope à la régularisation des déplacements. Chaque intégrale contient une norme discrète $|\cdot|$ définie par $|M| = \sqrt{\sum_{i,j} m_{i,j}^2}$, qui re-

présente une valeur absolue, une norme vectorielle ou une norme de Frobenius, selon le cas. Les paramètres λ et μ sont utilisés pour définir le compromis entre données et régularisations.

3.2 Énergie de reconstruction de la texture

La deuxième énergie E_T est une extension à la vidéo du travail de [1] utilisant directement les images couleur, et non canal par canal comme pour la reconstruction de la structure. Ici la recherche de patches optimaux n'est plus effectuée dans un voisinage spatial autour du défaut, mais dans un voisinage temporel (dans les images précédente u_b et suivante u_f). Alors que, dans [22] et [19], la recherche de patches 3D n'est pas limitée dans le temps, nous nous restreignons ici aux patches 2D dans les images avant et arrière. Dans l'image courante u , le pixel central x du patch $p_u(x)$ pour lequel on recherche le patch optimal dans les images voisines u_b et u_f se situe dans une zone \tilde{O} égale à l'extension de O d'une demi-largeur de patch, afin de propager les patches contenant suffisamment de données « saines » :

$$E_T(\Gamma_b, u, \Gamma_f) = \int_{\tilde{O}} \omega(x) \varepsilon [p_{u_b}(\Gamma_b(x)) - p_u(x)] dx + \int_{\tilde{O}} (1 - \omega(x)) \varepsilon [p_{u_f}(\Gamma_f(x)) - p_u(x)] dx \quad (3)$$

où Γ_b et Γ_f désignent les cartes de correspondances, respectivement, entre u et u_b et entre u et u_f , où $\omega(x) \in [0, 1]$ est un poids entre les deux reconstructions possibles de u

à partir des images avant ou arrière (voir Section 4.4) et où la fonction ε désigne la distance entre patches. En pratique, ε désigne une convolution entre la différence au carré des patches $(p_{u_b}(\Gamma_b(x)) - p_u(x))^2$ et un noyau gaussien g_a d'écart-type a pour une reconstruction non locale :

$$\varepsilon [p_{u_b}(\Gamma_b(x)) - p_u(x)] = \int_{\Omega_p} g_a(x_p) [u_b(\Gamma_b(x) - x_p) - u(x - x_p)]^2 dx_p \quad (4)$$

où $x_p \in \Omega_p$ correspond aux coordonnées d'un pixel dans le patch $p_u(x)$, relativement à son centre x . Minimiser la somme de E_S et E_T constitue un problème très complexe, sans preuve d'existence ou d'unicité d'une solution, mais on cherche seulement ici à obtenir une approximation numérique de la solution en minimisant E_S et E_T de manière alternée, en utilisant le résultat de la minimisation en u de l'un des deux termes pour initialiser la minimisation de l'autre terme.

4 Optimisation

Pour restaurer par *inpainting* des défauts de grande taille ou estimer le mouvement, une approche multirésolution est nécessaire, de la résolution la plus grossière ($L = L_{\max}$) vers la résolution la plus fine ($L = 0$). Dans le contexte de la vidéo, il est d'autant plus important de suivre cette stratégie pour pouvoir estimer correctement le mouvement.

4.1 Combinaison des deux modèles

Algorithme 1 - Reconstruction aux niveaux L et $L - 1$

- 1: **si** $L \geq L_{\min}^{\text{structure}}$ **alors**
 - 2: $v_b^*, u^*, v_f^* \leftarrow \operatorname{argmin}_{v_b, u, v_f} \{E_S(v_b, u, v_f)\}$
 - 3: $v_b, v_f \leftarrow \text{suréchantillonnage}(v_b^*, v_f^*)$
 - 4: $u \leftarrow u^*$
 - 5: **fin si**
 - 6: **si** $L \leq L_{\max}^{\text{texture}}$ **alors**
 - 7: $\Gamma_b^*, \Gamma_f^* \leftarrow \operatorname{argmin}_{\Gamma_b, \Gamma_f} \{E_T(\Gamma_b, u, \Gamma_f)\}$
 - 8: $\Gamma_b, u, \Gamma_f \leftarrow \text{suréchantillonnage}(\Gamma_b^*, u, \Gamma_f^*)$
 - 9: $u^* \leftarrow \operatorname{argmin}_u \{E_T(\Gamma_b, u, \Gamma_f)\}$
 - 10: $u \leftarrow u^*$
 - 11: **sinon**
 - 12: $u \leftarrow \text{suréchantillonnage}(u)$
 - 13: **fin si**
-

L'idée est de suffisamment sous-échantillonner les images pour pouvoir considérer que les mouvements dans la scène sont petits. À un niveau de résolution L proche de

L_{\max} , en utilisant un filtrage gaussien pour éliminer les hautes fréquences dans l'étape de sous-échantillonnage, la reconstruction de la structure fonctionne bien, alors que la reconstruction de la texture n'est pas efficace. En revanche, à plus haute résolution (L proche de 0), nous voulons mettre l'accent sur la texture. C'est pourquoi l'algorithme peut choisir un niveau maximal de résolution $L_{\max}^{\text{texture}}$ pour démarrer la reconstruction de la texture et un niveau minimal de résolution $L_{\min}^{\text{structure}}$ pour arrêter la reconstruction de la structure, avec $L_{\max}^{\text{texture}} \geq L_{\min}^{\text{structure}} - 1$ pour garantir qu'au moins une des reconstructions est effectuée à chaque résolution. Par conséquent, à chaque niveau de résolution, sauf au niveau de plus haute résolution $L = 0$, notre algorithme applique une seule reconstruction ou les deux à la fois.

Les différents suréchantillonnages sont opérés de façons différentes suivant les cas : une interpolation au plus proche voisin est appliquée aux cartes de correspondances, alors que l'interpolation est bicubique pour les autres variables. Les différentes minimisations de u , v_b , v_f , Γ_b et Γ_f de l'Algorithme 1 sont détaillées ci-après.

4.2 Estimation des mouvements

Afin de minimiser E_S par rapport au vecteur de mouvement v_b , on procède comme dans [23] en linéarisant l'argument des deux différences absolues dans (2). Cependant, cette linéarisation n'est possible que si les déplacements sont petits. C'est pourquoi un autre vecteur de mouvement constant v_b^0 (respectivement v_f^0) est introduit, proche de v_b , autour duquel ce dernier est estimé :

$$\begin{aligned} E_S(v_b, u, v_f) &= \mu \int_{\Omega} |\nabla u(x)| dx + \lambda \int_{\Omega} |Jv_b(x)| dx + \lambda \int_{\Omega} |Jv_f(x)| dx \\ &+ \int_{\Omega} |\nabla u_b(x + v_b^0(x)) \cdot [v_b - v_b^0](x) + u_b(x + v_b^0(x)) - u(x)| dx \\ &+ \int_{\Omega} |\nabla u_f(x + v_f^0(x)) \cdot [v_f - v_f^0](x) + u_f(x + v_f^0(x)) - u(x)| dx \end{aligned} \quad (5)$$

où $\nabla u_b(x + v_b^0(x)) \cdot [v_b - v_b^0](x) + u_b(x + v_b^0(x)) - u(x)$ sera noté $\rho(u, v_b, u_b, v_b)$ dorénavant (idem pour v_f). Minimiser E_S relativement au vecteur de mouvement v_b conduit à l'expression suivante :

$$\begin{aligned} v_b^* &= \operatorname{argmin}_{v_b} \max_y \int_{\Omega} |\rho(u, v_b, u_b, v_b)| dx \\ &+ \langle Jv_b | y \rangle - \iota_{B^\infty} \left(\frac{y}{\lambda} \right) \end{aligned} \quad (6)$$

où nous avons introduit la variable duale de v_b , $y : \Omega \rightarrow \mathbb{R}^{2 \times 2}$. Ce problème convexe peut être résolu par l'algorithme primal-dual de [10]. En remarquant que $Jv_b = [\nabla v_{b,1}, \nabla v_{b,2}]^\top$, l'opérateur adjoint de la jacobienne de v_b prend alors la forme $J^*y = -[\operatorname{div}([y_{1,1}, y_{1,2}]^\top), \operatorname{div}([y_{2,1}, y_{2,2}]^\top)]^\top$. Alors que l'opérateur proximal associé à y est une projection sur la boule L^∞ , l'opérateur proximal associé à v_b est un seuillage doux (voir [23] pour plus de détails) :

$$\begin{cases} y^{(n+1)} \leftarrow \operatorname{prox}_{\lambda \sigma \iota_{B^\infty}} (y^{(n)} + \sigma J \bar{v}_b^{(n)}) \\ v_b^{(n+1)} \leftarrow \operatorname{prox}_{\tau \rho(u, -, u_b)} (v_b^{(n)} - \tau J^* y^{(n+1)}) \\ \bar{v}_b^{(n+1)} \leftarrow v_b^{(n+1)} + \theta (v_b^{(n+1)} - v_b^{(n)}) \end{cases} \quad (7)$$

où $\sigma, \tau > 0$ sont des pas de temps et $\theta \in [0, 1]$. La minimisation de v_f s'effectue de la même manière.

4.3 Reconstruction de la structure

Une fois le mouvement estimé, le processus d'*inpainting* est obtenu grâce à la minimisation suivante :

$$\begin{aligned} u^* &= \operatorname{argmin}_u \int_{\Omega} |u_b(x + v_b(x)) - u(x)| dx \\ &+ \int_{\Omega} |u_f(x + v_f(x)) - u(x)| dx + \mu \int_{\Omega} |\nabla u(x)| dx \end{aligned} \quad (8)$$

Pour réécrire le problème convexe (8) avec des variables duales, la dépendance temporelle de u , v_b et v_f doit être clarifiée. En effet, en prenant 1 comme pas de temps entre deux images, et en introduisant la variable de temps t dans u , les fonctions u_b et u_f prennent la forme suivante :

$$\begin{aligned} u_b(x + v_b(x)) &= u(x + v_b(x, t), t - 1) = u(\varphi_b(x, t)) \\ u_f(x + v_f(x)) &= u(x + v_f(x, t), t + 1) = u(\varphi_f(x, t)) \end{aligned} \quad (9)$$

où φ_b et φ_f sont deux transformations supposées différentiables et inversibles, similaires aux cartes de correspondances Γ_b et Γ_f , en faisant l'hypothèse que $\varphi_b \circ \varphi_f = \varphi_f \circ \varphi_b = I_d$ presque partout. Avec ces nouvelles notations, (8) peut être réécrite comme :

$$\begin{aligned} u^* &= \operatorname{argmin}_u \max_z \left\langle \begin{array}{l} u \circ \varphi_b - u \\ u \circ \varphi_f - u \\ \nabla u \end{array} \middle| z \right\rangle \\ &- \iota_{B^\infty}(z_1) - \iota_{B^\infty}(z_2) - \iota_{B^\infty} \left(\frac{1}{\mu} [z_3, z_4]^\top \right) \end{aligned} \quad (10)$$

où nous introduisons la variable duale de u , $z : \Omega \rightarrow \mathbb{R}^4$. En notant K l'opérateur relatif à u dans le produit scalaire, cela conduit à l'opérateur adjoint K^* décrit dans [18] :

$$\begin{aligned} K^*z &= \det(J\varphi_f) z_1 \circ \varphi_f - z_1 + \det(J\varphi_b) z_2 \circ \varphi_b \\ &- z_2 - \operatorname{div}([z_3, z_4]^\top) \end{aligned} \quad (11)$$

La minimisation (10) peut également être effectuée en utilisant l'algorithme primal-dual de [10] :

$$\begin{cases} z_1^{(n+1)} \leftarrow \operatorname{prox}_{\sigma \iota_{B^\infty}} (z_1^{(n)} + \sigma (u \circ \varphi_b - \bar{u}^{(n)})) \\ z_2^{(n+1)} \leftarrow \operatorname{prox}_{\sigma \iota_{B^\infty}} (z_2^{(n)} + \sigma (u \circ \varphi_f - \bar{u}^{(n)})) \\ \begin{bmatrix} z_3^{(n+1)} \\ z_4^{(n+1)} \end{bmatrix} \leftarrow \operatorname{prox}_{\mu \sigma \iota_{B^\infty}} \left(\begin{bmatrix} z_3^{(n)} \\ z_4^{(n)} \end{bmatrix} + \sigma' \nabla \bar{u}^{(n)} \right) \\ u^{(n+1)} \leftarrow u^{(n)} - \tau K^* z^{(n+1)} \\ \bar{u}^{(n+1)} \leftarrow u^{(n+1)} + \theta (u^{(n+1)} - u^{(n)}) \end{cases} \quad (12)$$

où $\sigma, \sigma', \tau > 0$ sont des pas de temps et $\theta \in [0, 1]$. Pour minimiser une énergie similaire à (2), [7] répète les minimisations alternées pour l'*inpainting* de u et l'estimation d'un unique champ de déplacements v jusqu'à la convergence. A contrario, dans notre cas, les trois variables (v_b, v_f, u) sont minimisées simultanément, en appliquant des étapes proximales successives sur chaque variable. Les trois descentes suivant v_b, v_f et u donnent de meilleurs résultats que la première option, pour un temps de calcul équivalent.

4.4 Reconstruction de la texture avec estimation des cartes de correspondances

La minimisation de l'énergie de texture E_T se fonde sur les travaux de [1], où les cartes de correspondances Γ_b et Γ_f sont estimées en utilisant l'algorithme PatchMatch de [3] pour opérer une recherche non locale efficace sans être exhaustive, avec une distance L^2 entre patches. Par conséquent, pour Γ_b (respectivement Γ_f), on a, $\forall x \in \tilde{O}$:

$$\Gamma_b(x) = \operatorname{argmin}_{x_b \in \Omega} \int_{\Omega_p} g_a(x_p) [u_b(x_b - x_p) - u(x - x_p)]^2 dx_p \quad (13)$$

Pour chacune des images voisines u_b et u_f , chacun des termes du membre de droite de (3) peut être réécrit, sans prendre en compte le poids ω dans un premier temps, comme le cas extrême où l'on choisit seulement le patch le plus proche pour chaque pixel de la zone de défaut (voir [1] pour plus de détails), en introduisant une fonction de Dirac δ :

$$E_T^b(u, \Gamma_b) = \int_{\tilde{O}} \int_{\Omega} \delta(\Gamma_b(x) - x_b) \varepsilon[p_{u_b}(x_b) - p_u(x)] dx_b dx \quad (14)$$

Avec les changements de variables $x := x - x_p$ et $x_b := x_b - x_p$, qui opèrent deux translations, on tire de (4) et (14) :

$$E_T^b(u, \Gamma_b) = \int_{\tilde{O}} \int_{\Omega} m(x, x_b) [u_b(x_b) - u(x)]^2 dx_b dx \quad (15)$$

où :

$$m(x, x_b) = \int_{\Omega_p} g_a(x_p) \delta(\Gamma_b(x + x_p) - (x_b + x_p)) dx_p \quad (16)$$

dont l'intégrale sur Ω est égale à 1, puisque g_a est, par hypothèse, une gaussienne normalisée. En développant la différence au carré dans (15), u est aussi le minimiseur de l'énergie suivante, qui est égale à $E_T^b(u, \Gamma_b)$ à une constante près :

$$\tilde{E}_T^b(u, \Gamma_b) = \int_{\tilde{O}} \left[u(x) - \int_{\Omega} m(x, x_b) u_b(x_b) dx_b \right]^2 dx \quad (17)$$

ce qui nous fournit directement la solution des moyennes non locales définie, $\forall x \in \tilde{O}$, par :

$$\begin{aligned} u(x) &= \int_{\Omega} m(x, x_b) u_b(x_b) dx_b \\ &= \int_{\Omega_p} g_a(x_p) u_b(\Gamma_b(x + x_p) - x_p) dx_p \end{aligned} \quad (18)$$

Il s'agit ici du résultat pour une seule des deux images voisines. Pour prendre en compte les deux images, des moyennes pondérées peuvent être effectuées avec les poids $\omega(x)$ et $(1 - \omega(x))$. Moyenner les deux résultats ($\omega(x) = 0.5$, notée T_H) provoquera l'apparition de flou. À l'inverse, choisir le meilleur des deux résultats ($\omega(x) \in \{0, 1\}$, notée T_B) fera apparaître des artefacts spatiaux. Dans ce cas, le choix entre 0 et 1 pour les poids est effectué après avoir estimé les deux cartes de correspondances Γ_b et Γ_f . Ensuite, pour chaque pixel à reconstruire, le poids $\omega(x)$ est égal à 1 si la distance estimée $\varepsilon[p_{u_b}(\Gamma_b(x)) - p_u(x)]$ au meilleur patch dans l'image arrière est inférieure à la distance estimée $\varepsilon[p_{u_f}(\Gamma_f(x)) - p_u(x)]$ au meilleur patch dans l'image avant, et est égal à 0 sinon. Une solution intermédiaire consiste à appliquer une moyenne pondérée (notée T_W) en utilisant le ratio entre la distance estimée au meilleur patch dans l'image avant et la somme des deux distances estimées.

5 Expérimentations

Notre algorithme a été implémenté en Matlab et en C. Le niveau maximal de la pyramide multirésolution vaut $L_{\max} = 10$, et le facteur entre étages successifs est égal à $\sqrt{2}$. L'algorithme est d'abord testé sur trois images consécutives issues de séquences Middlebury, où un défaut a été introduit artificiellement dans l'image centrale (voir FIGURE 2). Il est ensuite testé en conditions réelles, sur des images numérisées issues de vieux films de la Cinéma-thèque de Toulouse.

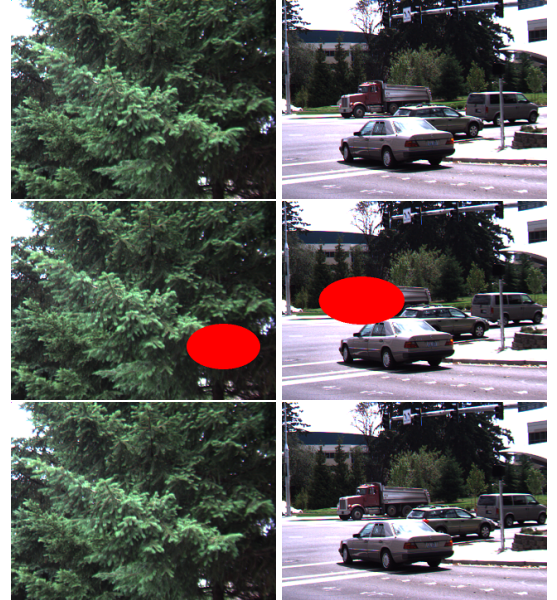


FIGURE 2 – Images 8 à 10 des séquences « Evergeen » (gauche) et « Dumptruck » (droite) de la base Middlebury. Le défaut est surligné en rouge dans l'image centrale.

Les vidéos des séquences complètes de Middlebury contenant 8 images et comportant le défaut artificiel, ainsi que les reconstructions correspondantes, sont disponibles en téléchargement sur le lien suivant : [ORASIS-VideoInpainting1](#).

5.1 Séquence « Evergreen »

Pour ce premier test, les résultats de la FIGURE 3 montrent que les méthodes où l'on ne reconstruit que la texture (3c, 3e, 3g) ne parviennent pas à restituer correctement les branches de l'arbre. L'approche avec le meilleur des deux patches (3c) provoque des artéfacts, alors que les approches utilisant un moyennage (3e et 3g) conduisent à un résultat où la texture est perdue au profit d'une homogénéité non désirée. La reconstruction de la structure (3b) est meilleure, mais encore un peu floue à cause de la diffusion, et certaines branches semblent un peu trop étirées. Enfin, la combinaison des deux reconstructions (3d, 3f, 3h) réussit à superposer une texture à la structure.

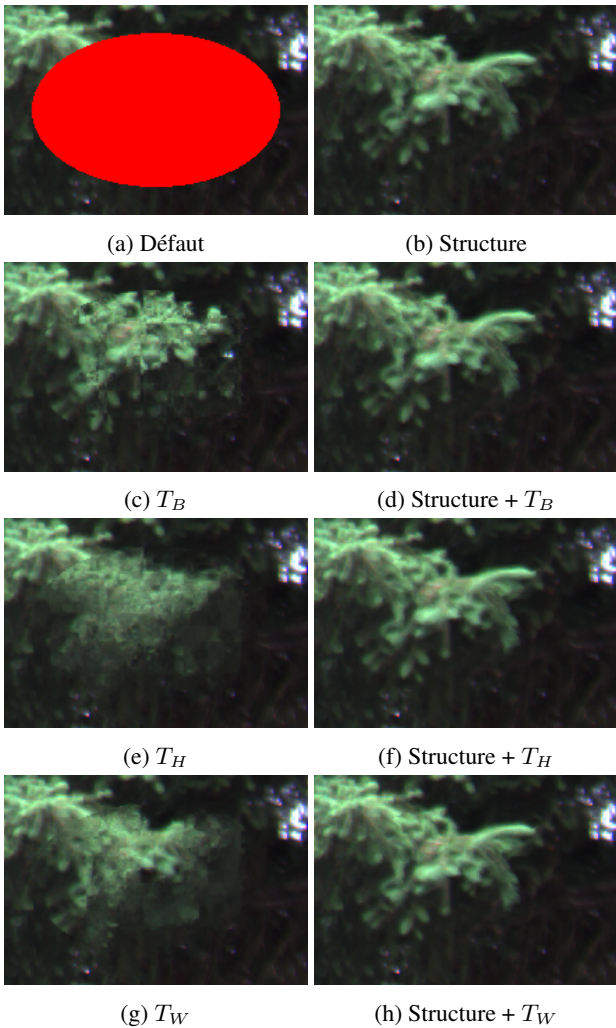


FIGURE 3 – Image centrale « Evergreen » : zooms sur les reconstructions de la zone défectueuse.

Ces premières interprétations visuelles sont confirmées par les mesures de qualité du TABLEAU 1. La structure est plus performante que la texture, tandis que notre algorithme donne les meilleurs résultats, avec un écart non négligeable. Par ailleurs, le choix de l'une quelconque des trois reconstructions de la texture, associée à la reconstruction de la structure, n'a pas de réel impact sur le résultat.

	Sans structure	Structure
Sans texture	18.66 - 0.943	37.24 - 0.991
Texture T_B	30.94 - 0.978	40.37 - 0.994
Texture T_H	33.37 - 0.980	40.47 - 0.994
Texture T_W	33.56 - 0.982	40.50 - 0.994

TABLEAU 1 – Mesures de PSNR - SSIM pour les différentes reconstructions de l'image centrale « Evergreen ».

5.2 Séquence « Dumptruck »

Pour ce deuxième test, les résultats de la FIGURE 4 montrent que la reconstruction de la structure (4b) entraîne une déformation du camion, qui semble osciller verticalement dans la vidéo. De plus, il reste une trace derrière la voiture se dirigeant vers la droite. Les reconstructions de la texture à l'aide de moyennes (4e et 4g) conduisent à une reconstruction floue du camion. Même en utilisant le meilleur patch des deux images adjacentes (4c), dont le résultat semble être bon sur l'image reconstruite seule, la vidéo montre que l'algorithme choisit l'image la plus proche, en terme de distance entre patches, et reste « verrouillé ». C'est pourquoi il semble y avoir un manque de mouvement dans la vidéo, à l'intérieur de la zone de défaut. Notre algorithme (4d, 4f, 4h) fournit de meilleurs résultats (il reste quand même une trace à l'arrière de la voiture), grâce à l'apport de la texture qui permet de conserver un bon flux optique.

D'autre part, les mesures de qualité du TABLEAU 2 montrent que notre algorithme améliore significativement les résultats, en comparaison des reconstructions obtenues avec la structure seule ou avec la texture seule.

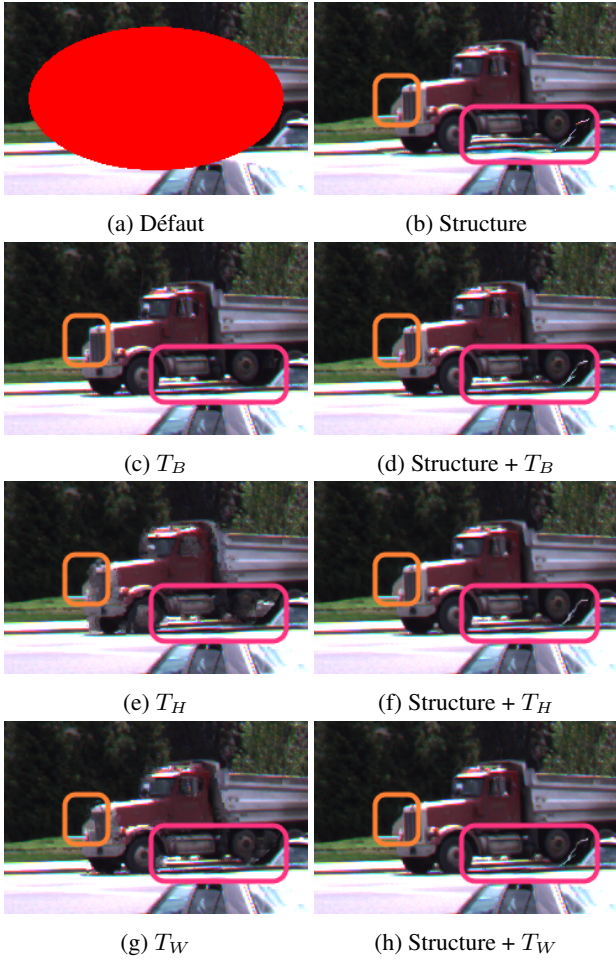


FIGURE 4 – Image centrale « Dumpruck » : zooms sur les reconstructions de la zone défectueuse. Dans la zone encadrée en magenta, la structure seule (4b) déforme le dessous du camion. La reconstruction près de la voiture de droite reste problématique dans tous les cas. Dans la zone encadrée en orange, la texture avec choix du meilleur patch (4c) décale le camion vers la gauche.

	Sans structure	Structure
Sans texture	17.53 - 0.937	27.76 - 0.975
Texture T_B	33.16 - 0.987	38.11 - 0.994
Texture T_H	34.35 - 0.989	38.87 - 0.994
Texture T_W	34.31 - 0.988	38.58 - 0.994

TABLEAU 2 – Mesures de PSNR - SSIM pour les différentes reconstructions de l'image centrale « Dumpruck ».

5.3 Séquence de la Cinémathèque



FIGURE 5 – Sélection manuelle d'un masque englobant les zones défectueuses : image comportant des traces blanches (gauche), et sélection en rouge du masque englobant les défauts (droite).

Les tests sur des séquences comportant un défaut artificiel ont montré que la combinaison de la structure et de la texture améliorerait effectivement les reconstructions. L'enjeu de l'algorithme proposé est de pouvoir corriger les défauts de vidéos ayant subi une réelle dégradation. Nous utilisons pour cela les images numérisées de vieux films argentiques provenant du fonds de la Cinémathèque de Toulouse (voir FIGURE 1). La première étape consiste à déterminer manuellement les zones défectueuses sous la forme d'un masque binaire (voir FIGURE 5).

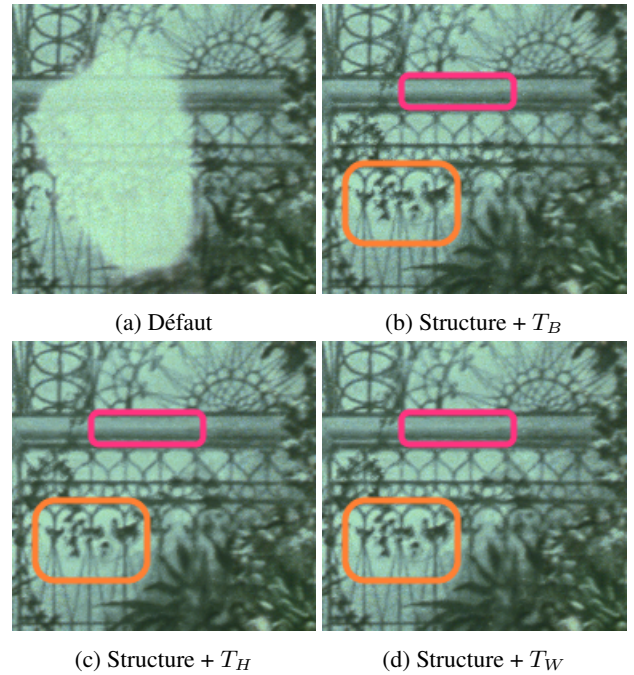


FIGURE 6 – Image de la FIGURE 5 : zooms sur trois reconstructions de l'une des zones détériorées. Dans la zone encadrée en magenta, l'utilisation de la moyenne des deux patches (6c) provoque un lissage excessif. A contrario, dans la zone encadrée en orange, le choix du meilleur des deux patches (6b) a pour effet de trop accentuer les contours.

Une fois le défaut détecté, les trois variantes de notre algorithme, appliquées à la séquence contenant l'image de la FIGURE 5, fournissent trois reconstructions de qualité globalement satisfaisante (voir FIGURE 6). Malgré tout, la question d'un surplus de flou (6c) ou de la présence d'artéfacts (6b) se pose encore, en particulier vis-à-vis des professionnels du cinéma, pour qui la préservation du grain argentique semble cruciale. La solution intermédiaire (6d) peut donc constituer un bon compromis, car elle lisse légèrement les formes tout en évitant les artéfacts et en préservant suffisamment le grain. Sur les vidéos, la différence entre ces trois reconstructions est difficile à percevoir avec un défaut présent sur une seule image (laps de temps trop court). Les différentes vidéos sont disponibles en téléchargement sur le lien suivant : [ORASIS-VideoInpainting2](#).

6 Conclusion et perspectives

Dans cet article, nous avons montré que la combinaison d'une reconstruction de la structure par des approches de diffusion et d'une reconstruction de la texture permettait d'améliorer les résultats de l'*inpainting* vidéo, en termes de qualité visuelle et de métriques. Pour aller plus loin, il serait intéressant d'affiner le modèle de flux optique en utilisant la variation totale généralisée, comme dans [20]. D'autres reconstructions non locales de la texture peuvent également être envisagées, par exemple avec le filtre médian, ou en utilisant des gradients de patches ou de nouvelles régularisations de patches. D'autre part, nous envisageons de recourir au *deep learning*, soit en amont pour la détection des défauts, soit pour le calcul des opérateurs proximaux de la reconstruction elle-même.

Références

- [1] P. Arias, G. Facciolo, V. Caselles, and G. Sapiro. A Variational Framework for Exemplar-based Image Inpainting. *IJCV*, 93, 2011.
- [2] G. Aubert, R. Deriche, and P. Kornprobst. Computing Optical Flow via Variational Techniques. *SIAM Journal on Applied Mathematics*, 60, 1999.
- [3] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman. PatchMatch: A Randomized Correspondence Algorithm for Structural Image Editing. *ACM Transactions on Graphics*, 28, 2009.
- [4] M. Bertalmio, G. Sapiro, V. Caselles, and C. Ballester. Image Inpainting. In *Proceedings of SIGGRAPH*, 2000.
- [5] A. Buades and J.-L. Lisani. Video Denoising with Optical Flow Estimation. *IPOL*, 8, 2018.
- [6] A. Bugeau and M. Bertalmio. Combining Texture Synthesis and Diffusion for Image Inpainting. In *Proceedings of VISAPP*, 2009.
- [7] M. Burger, H. Dirks, and C.-B. Schönlieb. A Variational Model for Joint Motion Estimation and Image Reconstruction. *SIAM Journal on Imaging Sciences*, 11, 2018.
- [8] P. Buysens, M. Daisy, D. Tschumperlé, and O. Lézoray. Exemplar-based Inpainting: Technical Review and New Heuristics for Better Geometric Reconstructions. *IEEE Transactions on Image Processing*, 24, 2015.
- [9] F. Cao, Y. Gousseau, S. Masnou, and P. Pérez. Geometrically guided exemplar-based inpainting. *SIAM Journal on Applied Mathematics*, 4, 2011.
- [10] A. Chambolle and T. Pock. A First-order Primal-dual Algorithm for Convex Problems with Applications to Imaging. *JMIV*, 40, 2011.
- [11] T. F. Chan and J. Shen. Local Inpainting Models and TV Inpainting. *SIAM Journal on Applied Mathematics*, 62, 2001.
- [12] J.-P. Cocquerez, L. Chanas, and J. Blanc-Talon. Simultaneous Inpainting and Motion Estimation of Highly Degraded Video-sequences. In *Proceedings of SCIA*, 2003.
- [13] A. Criminisi, P. Pérez, and K. Toyama. Region Filling and Object Removal by Exemplar-Based Image Inpainting. *IEEE Transactions on Image Processing*, 13, 2004.
- [14] V. Do, G. Lebrun, L. Malapert, C. Smet, and D. Tschumperlé. Inpainting d'images couleurs par lissage anisotrope et synthèse de textures. In *Proceedings of RFIA*, 2006.
- [15] A. A. Efros and T. K. Leung. Texture Synthesis by Non-parametric Sampling. In *Proceedings of ICCV*, 1999.
- [16] B. K. Horn and B. G. Schunck. Determining Optical Flow. *Artificial Intelligence*, 17, 1981.
- [17] F. Lauze and M. Nielsen. A Variational Algorithm For Motion Compensated Inpainting. In *Proceedings of BMVC*, 2004.
- [18] F. Lauze and M. Nielsen. On Variational Methods for Motion Compensated Inpainting. *arXiv preprint*, 2009.
- [19] T. Le, A. Almansa, Y. Gousseau, and S. Masnou. Motion-Consistent Video Inpainting. In *Proceedings of ICIP*, 2017.
- [20] R. March and G. Riey. Analysis of a Variational Model for Motion Compensated Inpainting. *Inverse Problems & Imaging*, 11, 2017.
- [21] S. Masnou and J.-M. Morel. Level Lines Based Disocclusion. In *Proceedings of ICIP*, 1998.
- [22] A. Newson, A. Almansa, M. Fradet, Y. Gousseau, and P. Pérez. Video Inpainting of Complex Scenes. *SIAM Journal on Imaging Sciences*, 7, 2014.
- [23] C. Zach, T. Pock, and H. Bischof. A Duality Based Approach for Realtime TV-L1 Optical Flow. In *Proceedings of Joint Pattern Recognition Symposium*, 2007.