

Mécanisme de synchronisation en tatouage audio pour des perturbations désynchronisantes à forte dérive

Cléo BARAS, Nicolas MOREAU, Bassem ZAYEN

GET - Télécom Paris, Département TSI
46 rue Barrault, 75634 Paris Cédex 13, France
{baras, moreau, zayen}@tsi.enst.fr

Résumé – Pour satisfaire les enjeux de robustesse des systèmes de tatouage, nous proposons dans cet article un mécanisme de synchronisation quasi temps-réel. Ce mécanisme tolère toutes opérations désynchronisantes, à condition qu’elles soient modélisables par un rééchantillonnage du signal audio tatoué. Ces opérations induisent une modification de l’échelle des temps entre l’émission et la réception et peuvent être caractérisées par un facteur de dérive, écart relatif entre les fréquences d’échantillonnage. Des résultats expérimentaux sur des signaux réels montrent que le mécanisme de synchronisation proposé maintient la fiabilité de transmission obtenue à un débit donné à une valeur quasi-constante indépendante de la dérive. Une transmission de fiabilité n’excédant pas 10^{-3} peut donc être obtenue pour des débits inférieurs à 50 bits/s y compris en présence de dérives allant jusqu’à 2%.

Abstract – In this paper, we propose a real-time synchronization mechanism for audio watermarking systems. This mechanism pursues the goal of designing a watermarking system which is robust to all desynchronizing operations, provided that these operations can be modeled as a resampling of the watermarked audio signal. Time is rescaled between the embedder and the receiver with respect to a scaling drift, defined as the relative distance between the sampling rates. Experimental results on real audio signals show that the transmission reliability obtained for a given transmission rate remains steady and does not depend on the scaling factor. Consequently, a real-time transmission with 10^{-3} reliability is feasible for transmission rate lower than 50 bps, even when the scaling drift reaches 2%.

1 Introduction

Si les techniques de tatouage sont apparues pour répondre au problème de protection des droits d’auteur, elles sont aujourd’hui aussi utilisées pour étendre le contenu du signal audio, en y insérant une information supplémentaire à destination d’une application cible. Cette application cible est par exemple l’animation d’un clone virtuel pour le projet RNRT¹ ARTUS² [1]. Dans ce contexte applicatif particulier, les performances d’un système de tatouage audio dépendent de quatre critères : la transparence du tatouage, le couple débit / fiabilité de transmission, la robustesse aux perturbations licites subies les signaux audio et le coût en temps de calcul des traitements. Parmi l’ensemble des perturbations envisageables, les plus dégradantes sont sans doute les perturbations dites désynchronisantes, auxquelles nous nous intéressons dans ce papier.

Plusieurs types de désynchronisation peuvent être répertoriés, allant de l’introduction d’un simple retard lors d’un filtrage à la dilation ou la contraction de la durée du signal audio lors de Conversions Analogique-Numérique (CAN) ou d’opération de "time stretching" [2]. Ces perturbations peuvent être modélisées par un rééchantillonnage du signal audio tatoué. Ce rééchantillonnage introduit une dérive entre l’échelle des temps à l’émission et à la réception, modifiant de fait de la localisation des bits transmis et leur durée. Cette dérive, définie comme l’écart relatif entre les fréquences d’échantillonnage, peut prendre des valeurs faibles - typiquement de l’ordre de 0,1% si l’on veut utiliser le canal audio pour transmettre de l’information entre deux PCs par exemple - ou nettement plus

élevée - de l’ordre de quelques pour cent dans des applications de radiodiffusion. En effet, une émission radio peut être entrecoupée de séquences sonores, préalablement tatouées et devant être diffusées dans un temps imparti. Le "time stretching" permet alors de modifier pour l’occasion la durée de la séquence de sorte à satisfaire la contrainte de temps de diffusion.

L’état de l’art en matière de synchronisation des systèmes de tatouage propose diverses méthodes que l’on peut regrouper en deux classes. La première exploite des formes d’onde étalées spectralement pour estimer la désynchronisation introduite (différence d’échantillonnage entre l’émission et la réception [3] en audio ou distorsion géométrique en image). Souvent adaptées aux opérations de désynchronisation illicites, elle nécessite un temps de calcul relativement important, qui ne permet pas d’envisager un traitement temps-réel du récepteur. La seconde classe consiste à extraire du signal des points d’intérêt [4, 5], qui définissent une cartographie du signal. L’extraction de cette "carte" à la réception permet là encore d’estimer et d’inverser l’opération désynchronisante pour procéder ensuite à la détection de l’information cachée. Malheureusement, ces points d’intérêt étant souvent choisis en fonction des caractéristiques du signal, les systèmes de tatouage qui en découlent induisent un débit de transmission variable.

Dans cet article, nous proposons un mécanisme de synchronisation quasi-temps réel et induisant un débit de transmission fixe. Ce mécanisme est dédié aux systèmes de tatouage dans le domaine temporel conçu dans un contexte de transmission d’information. Il vise à garantir la robustesse de la transmission en présence de perturbations désynchronisantes dont la dérive peut prendre des valeurs importantes mais ne présente pas de variations. Nous appliquerons plus particulièrement ce

¹Réseau National de Recherche en Télécommunication

²Animation Réaliste par Tatouage à Usage des Sourds

mécanisme à une stratégie de tatouage de type additif par étalement de spectre décrite dans [6], dont nous rappellerons les principes dans la première partie. Le mécanisme de synchronisation fera l'objet de la seconde partie. Nous présenterons pour finir des résultats expérimentaux sur des signaux réels qui permettront de conclure quant aux performances du mécanisme de synchronisation.

2 Système de tatouage

Le système de tatouage considéré est le système de tatouage de type additif par étalement de spectre proposé par Larbi dans [6], dont le schéma de principe est rappelé figure 1.

L'information à émettre est supposée être une séquence de L symboles $\{s_l\}_{l=0..L-1}$ à valeur dans $\{0..M-1\}$. Son insertion requière un dictionnaire d'émission $\mathcal{D} = \{\mathbf{d}_m\}_{m=0..M-1}$, contenant M vecteurs blancs orthogonaux de durée finie N_s , de puissance unité et limités à la bande de fréquence $[0; F_c]$. Sur chaque temps symbole $[lN_s; (l+1)N_s - 1]$, où doit être inséré l'information s_l , le signal modulé \mathbf{v} est choisi égal au vecteur \mathbf{d}_{s_l} . Il est ensuite mis en forme spectralement par le filtre $H(f)$ pour obtenir le signal de tatouage \mathbf{t} . Ce filtre est conçu de sorte que la Densité Spectrale de Puissance (DSP) de \mathbf{t} coïncide avec un seuil de masquage, limite fréquentielle caractérisant la contrainte d'inaudibilité du signal \mathbf{t} en présence de \mathbf{x} . Ce seuil est calculé en appliquant un modèle psychoacoustique (dérivé du modèle MPEG n°1) au signal audio. Le signal audio tatoué \mathbf{y} est finalement obtenu en ajoutant le signal audio \mathbf{x} et le signal de tatouage \mathbf{t} .

La réception est basée sur une procédure d'égalisation qui consiste à inverser la mise en forme spectrale de l'émission puis à réestimer le signal modulé à l'aide d'un filtre de Wiener. Le signal audio tatoué reçu $\hat{\mathbf{y}}$ (après d'éventuelles perturbations dans le canal) est donc filtré par $\frac{1}{\hat{H}(f)}$, où $\hat{H}(f)$ est le filtre de mise en forme calculé à partir de $\hat{\mathbf{y}}$. Un filtre de Wiener est ensuite construit de sorte à minimiser l'erreur quadratique entre le signal estimé $\hat{\mathbf{v}}$ et le signal modulé \mathbf{v} relativement aux coefficients du filtre. Ces coefficients \mathbf{w} sont donnés par : $\mathbf{w} = R_{\hat{\mathbf{z}}}^{-1} r_{\mathbf{v}}$, où $R_{\hat{\mathbf{z}}}$ et $r_{\mathbf{v}}$ sont respectivement la matrice et la fonction d'autocovariance du signal $\hat{\mathbf{z}}$ et du signal \mathbf{v} à estimer.

La décision quant à la séquence binaire reçue est effectuée à l'aide d'un démodulateur par corrélation qui sélectionne le vecteur d'un dictionnaire de réception $\hat{\mathcal{D}} = \{\hat{\mathbf{d}}_m\}_{m=0..M-1}$ dont la corrélation avec le signal estimé $\hat{\mathbf{v}}$ est maximum. Ce dictionnaire doit refléter l'ensemble des distorsions introduites par le canal sur les vecteurs du dictionnaire. Le signal n'ayant pour l'instant subi aucune désynchronisation, $\hat{\mathcal{D}}$ est donc identique au dictionnaire d'émission \mathcal{D} .

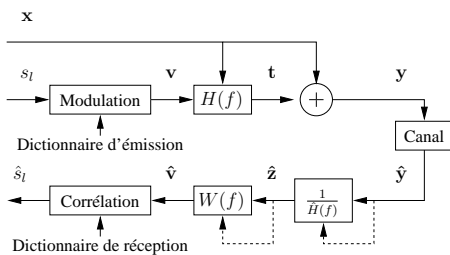


FIG. 1 – Schéma de principe du système de tatouage

3 Mécanisme de synchronisation

3.1 Modélisation de la désynchronisation

Les opérations désynchronisantes du canal, telles que les opérations de contraction-dilatation, peuvent être modélisées par l'ajout d'un retard et le rééchantillonnage du signal audio tatoué \mathbf{y} (initialement échantillonné à la fréquence F_e) à la fréquence \hat{F}_e . En exploitant la relation de conversion d'un signal numérique en signal analogique, le signal audio tatoué rééchantillonné $\hat{\mathbf{y}}$ peut être exprimé sous la forme :

$$\hat{y}(\hat{n} + \hat{n}_0) = \sum_{q=-\infty}^{\infty} y(n(\hat{n}) + q) \sin_c(q + \delta(\hat{n})),$$

où \hat{n}_0 est le retard à l'origine, $n(\hat{n}) = \left\lfloor \frac{F_e}{\hat{F}_e} \hat{n} \right\rfloor$ ($\lfloor x \rfloor$ dénotant l'arrondi de x à l'entier le plus proche) modélise la transformation subie par l'échelle des temps entre l'émetteur et le récepteur et $\delta(\hat{n}) = n(\hat{n}) - \left\lfloor \frac{F_e}{\hat{F}_e} \hat{n} \right\rfloor$ est l'erreur liée à l'arrondi.

Ces opérations modifient donc l'amplitude des échantillons du signal audio tatoué mais surtout la location de chaque symbole et la durée pendant laquelle ils sont insérés. Le mécanisme de synchronisation, pour être efficace, doit donc restituer aussi bien le rythme que les formes prises par l'information tatouée, traduites par un dictionnaire de réception reflétant les perturbations introduites par la désynchronisation. Ces étapes requièrent finalement l'estimation de la fréquence \hat{F}_e .

3.2 Principe du mécanisme

Le mécanisme de synchronisation que nous proposons, représenté figure 2, peut être schématisé sous la forme de 3 étapes :

- une *phase d'initialisation* servant à l'estimation de la fréquence \hat{F}_e et du retard \hat{n}_0 . De cette phase est déduite la dérive $d = \frac{\hat{F}_e - F_e}{F_e}$.
- une *phase d'acquisition*, permettant la synchronisation du récepteur sur des points de référence, cadencant le message binaire émis et l'actualisation de la valeur de \hat{F}_e .
- une *phase de poursuite*, réalisant la synchronisation fine de la localisation des symboles et leurs détections avec un dictionnaire de réception adapté à la désynchronisation.

3.2.1 Principe général

Cette procédure de synchronisation nécessite de modifier la structure du signal de tatouage en ajoutant différents patterns de synchronisation comme schématisé à la figure 3. Un pattern de synchronisation \mathbf{p}_i de durée N_i est insérée M_i fois en amont du signal de tatouage à intervalle de $N_i + N_z$ échantillons. Le signal de tatouage est ensuite découpé sous la forme de *messages* de longueur N_m , contenant un pattern de synchronisation \mathbf{p}_a de durée N_a servant de points de référence et une partie de la séquence binaire à tatouer, de L_m symboles, appelée *sous-séquence*. Les patterns sont choisies étalées spectralement dans la bande de fréquence $[0, F_c]$ et de puissance unité.

La stratégie de tatouage utilisée pour leurs insertions dans le signal audio est identique à celle du système de tatouage décrit section 2 : ces patterns sont filtrées par le filtre de mise en forme psychoacoustique $H(f)$. Elles seront estimées à la

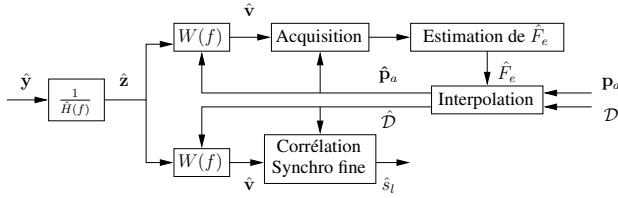


FIG. 2 – Mécanisme de synchronisation pour les phases d'acquisition et de poursuite

réception pour les besoins du mécanisme de synchronisation en utilisant la procédure d'égalisation préalablement décrite : le filtre d'inversion du canal $1/\hat{H}(f)$ (dont résulte le signal \hat{z}) puis un filtre de Wiener $W(f)$. Le filtre $\hat{H}(f)$, obtenu suivant un modèle psychoacoustique appliqué au signal audio tatoué \hat{y} , est peu sensible à la désynchronisation (puisqu'il est homogène à une DSP). Par contre, le filtre de Wiener $W(f)$ dépend de la pattern à estimer et donc de la modification qu'elle a pu subir pendant l'opération désynchronisante : il doit donc être calculé spécifiquement en fonction de l'interpolation à la fréquence \hat{F}_e de la pattern à estimer. Plus généralement, on dira par la suite que $W(f)$ est calculé spécifiquement pour estimer le signal \mathbf{v} si ces coefficients choisis de sorte à minimiser l'erreur quadratique entre le signal résultant du filtrage $\hat{\mathbf{v}}$ et le signal à estimer \mathbf{v} . Ces coefficients seront : $\mathbf{w} = R_{\hat{\mathbf{z}}}^{-1} r_{\mathbf{v}}$, où $r_{\mathbf{v}}$ est la fonction d'autocovariance de \mathbf{v} .

3.2.2 Phase d'initialisation

La phase d'initialisation consiste d'abord à estimer la fréquence de réception \hat{F}_e à l'aide des M_i patterns de synchronisation insérées en amont du signal de tatouage. Pour des raisons d'implémentation, ces patterns sont pour l'instant recherchés sur une durée limitée N_r au début du signal audio.

\hat{F}_e n'ayant pas encore été estimée, $W(f)$ est spécifiquement calculé pour la pattern \mathbf{p}_i . P_i versions interpolées de \mathbf{p}_i , notées $\hat{\mathbf{p}}_i$, sont évaluées à des fréquences \hat{F}_i , couvrant l'intervalle des fréquences de réception tolérées par le système de synchronisation. Pour chaque version interpolée $\hat{\mathbf{p}}_i$, la fonction d'intercovariance entre le signal $\hat{\mathbf{v}}$ (issu du filtrage de Wiener) et $\hat{\mathbf{p}}_i$ est calculée. De l'ensemble de ces P_i fonctions d'intercovariance sont ensuite seulement extraites les M_i valeurs maximales des corrélations. La durée entre ces M_i pics comparée à la durée imposée à l'émission entre les patterns indique la première estimation de \hat{F}_e . Cette estimation est ensuite affinée en traitant le signal audio une seconde fois : $W(f)$ est maintenant calculé spécifiquement pour la version interpolée de \mathbf{p}_i à la fréquence \hat{F}_e . La fonction d'intercorrélation entre la pattern interpolée et le signal $\hat{\mathbf{v}}$ résultant du filtrage est calculée, mettant en évidence M_i pics de corrélation. Un calcul de la durée entre les pics affine l'estimation de \hat{F}_e .

La phase d'initialisation se termine par l'estimation du retard à l'origine \hat{n}_0 . Ce retard est évalué en détectant la pattern de synchronisation \mathbf{p}_a marquant le début du premier message. Connaissant l'estimation de \hat{F}_e , une version interpolée de \mathbf{p}_a , notée $\hat{\mathbf{p}}_a$, est calculée et utilisée pour choisir les coefficients spécifiques de $W(f)$. Une méthode de corrélation par fenêtre glissante entre le signal $\hat{\mathbf{v}}$ résultant du filtrage est utilisée pour localiser l'instant de début de la pattern, noté \hat{n}_1 , sur lequel se synchronise le récepteur.

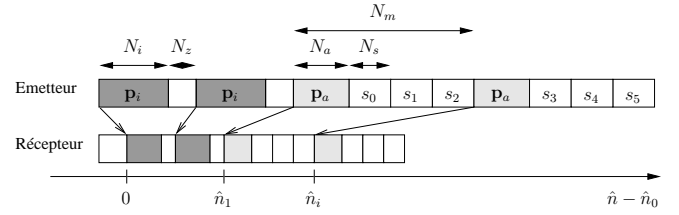


FIG. 3 – Structure du signal modulé (localisant les patterns de synchronisation) à l'émission et à la réception

3.2.3 Phase d'acquisition

La phase d'acquisition consiste à détecter les patterns de synchronisation \mathbf{p}_a insérées régulièrement dans le signal de tatouage pour synchroniser le récepteur. Ces patterns servent en effet de points de référence marquant le début de chaque message. La pattern du premier message ayant été détectée durant la phase d'initialisation, cette étape n'intervient donc qu'à partir de la détection du deuxième message.

Connaissant la position de la pattern du i -ème message, notée \hat{n}_i , et une estimation de la fréquence \hat{F}_e , la pattern de synchronisation précédant le $i + 1$ -ème message est recherchée autour de sa position théorique : $\hat{n}_{i+1} = \hat{n}_i \left[\frac{\hat{F}_e}{F_e} N_m \right]$, évitant ainsi une recherche exhaustive sur toute la longueur du message. Le filtre de Wiener est recalculé pour s'adapter à la version interpolée de \mathbf{p}_a par \hat{F}_e . La détection du pic de la fonction d'intercorrélation entre le signal filtré résultant $\hat{\mathbf{v}}$ et la pattern interpolée de \mathbf{p}_a par \hat{F}_e indique l'instant effectif de début du $i + 1$ -ème message.

La valeur estimée \hat{F}_e de la fréquence d'échantillonnage de réception est ensuite actualisée en évaluant la durée entre la pattern de synchronisation détectée et celle du premier message : $\hat{F}_e = \frac{\hat{n}_{i+1} - \hat{n}_1}{(i+1)N_m} F_e$.

3.2.4 Phase de poursuite

La pattern de synchronisation en amont du i -ème message ayant été localisée à l'instant \hat{n}_i , il est alors possible de procéder à la détection de la sous-séquence binaire qui lui fait suite. Connaissant l'estimation de \hat{F}_e , nous pouvons déduire les paramètres du processus de détection :

- le dictionnaire de réception adapté $\hat{\mathcal{D}}$ reflétant les distortions introduites par la désynchronisation : ce dictionnaire contient l'ensemble des vecteurs du dictionnaire d'émission \mathcal{D} interpolés à la fréquence estimée \hat{F}_e ,
- le filtre de Wiener spécifiquement adapté à l'estimation du signal modulé ayant subi l'interpolation par \hat{F}_e ,
- la durée d'insertion des symboles à la réception $\hat{N}_s = \frac{\hat{F}_e}{F_e} N_s$ ainsi que leur localisation théorique : le l -ième symbole du i -ème message débute à l'échantillon $\hat{n}_l = \hat{n}_i + \left[\frac{\hat{F}_e}{F_e} (N_a + lN_s) \right]$.

Une erreur sur l'estimation de \hat{F}_e ou une imprécision liée au calcul de l'arrondi pouvant néanmoins modifier légèrement la localisation du début de chaque symbole, on procède donc à une synchronisation fine de type "early-late gate" [7] pour préciser la localisation de chaque symbole reçu et décider conjointement de sa valeur. Cette synchronisation fine consiste à évaluer la métrique de décision sur un voisinage de la localisation théorique du symbole de N_e échantillons. La valeur maximale

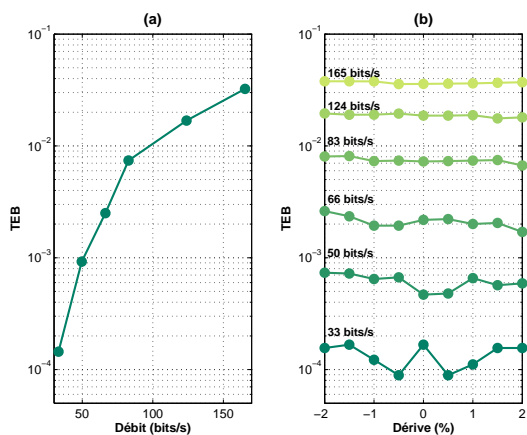


FIG. 4 – TEB obtenu pour différents débits par le système de tatouage avec le mécanisme de synchronisation proposé pour différentes configurations de canal : (a) sans perturbation, (b) opérations d'interpolation avec différentes valeurs de dérive.

de la corrélation indique à la fois la localisation précise du début du symbole reçu et sa valeur.

4 Evaluation des performances

Les performances du mécanisme de synchronisation sont jugées en évaluant la fiabilité de transmission et le coût en temps de calcul du système de tatouage, implémenté en C, face à des opérations de compression-dilatation pour différentes valeurs de la dérive. L'évaluation de la fiabilité consiste à mesurer le taux d'erreur binaire (TEB) obtenu par tatouage d'un ensemble de signaux tests par une séquence binaire de $L = 100000$ bits. Cet échantillon test est constitué de 20 signaux audio de style divers, échantillonnés à $F_e = 44.1$ kHz. La mesure du TEB fournit donc une estimation de la probabilité d'erreur de transmission à une précision de $5 \cdot 10^{-4}$ pour un taux de confiance de 70%. Le coût en temps de calcul est quant à lui évalué comme le rapport entre le temps nécessaire à la détection du tatouage et la durée du signal traité (en secondes).

Le dictionnaire utilisé pour le système de tatouage contient $M = 4$ vecteurs, étalés jusqu'à $F_c = 6$ kHz. Il permet d'obtenir un tatouage légèrement audible mais n'induisant pas de gêne auditive lors de l'écoute du signal tatoué (note entre 0 et -1 sur l'échelle perceptuelle définie par la recommandation UIT-R BS 562). Le débit de transmission utile est : $R = \frac{\log_2(M)L_m N_s}{N_m} F_e$.

Les paramètres du mécanisme de synchronisation sont les suivants : les durées de patterns sont $N_i = 12N_{hw}$, $N_z = 2N_{hw}$, $N_a = 15N_{hw}$, où $N_{hw} = 512$ est la taille des fenêtres utilisées pour le calcul des filtres $H(f)$ et $W(f)$, leur taux de répétition est $M_i = 3$, la durée des messages est $N_m = F_e$, la durée de recherche pour la phase d'initialisation est $N_r = F_e$, celle de la localisation de chaque symbole est $N_e = 3$ et les paramètres pour l'interpolation de la pattern de synchronisation pour la phase d'initialisation sont : $P_i = 21$ et $\tilde{F}_i = F_e \pm 100i$.

Un rappel des performances du système en terme de débit de transmission et de TEB dans le cas d'un canal sans perturbation est proposé figure 4(a). La figure 4(b) permet d'évaluer les performances du système pour différents débits face à une opération de désynchronisation (de type interpolation) en présentant les TEB en fonction de la dérive. Le mécanisme de synchro-

nisation fait ainsi état de son efficacité puisque les TEB obtenus varient peu en fonction de la dérive et sont d'un ordre de grandeur similaire à ceux obtenus pour une transmission sans désynchronisation.

Le coût lié au temps de calcul est évalué à 0.5 fois le temps réel sans prendre en compte la phase d'initialisation. Cette phase, la plus coûteuse, n'est par contre effectuée qu'une unique fois en début de transmission : elle ne participe pas au "régime courant" de la détection, constitué des phases d'acquisition et de poursuite. Par la suite, elle devra être adaptée au besoin d'une réelle application temps-réel (puisque la recherche du début du tatouage ne doit pas être limitée à quelques secondes).

5 Conclusions et perspectives

Dans cet article, nous avons présenté un mécanisme de synchronisation quasi-temps réel adapté aux systèmes de tatouage dans le domaine temporel dédiés à la transmission d'information. Ce mécanisme est conçu pour garantir la robustesse de la transmission à toutes opérations de désynchronisation, à condition qu'elles puissent être modélisées par un rééchantillonnage du signal audio tatoué. Il est basé sur l'ajout de patterns de synchronisation, permettant d'estimer la dérive d'échantillonnage, même forte, pour adapter les paramètres du récepteur. Des résultats expérimentaux sur des signaux réels ont montré que la fiabilité de transmission obtenue pour un canal sans perturbation est correctement maintenue si le canal devient une désynchronisation à forte dérive, attestant ainsi de la performance du mécanisme proposé. Ce mécanisme étant susceptible d'être utilisé avec tous systèmes de tatouage réalisant une insertion dans le domaine temporel, des études sont actuellement menées pour évaluer ces performances pour différents types de récepteur (par filtre blanchissant par exemple) et pour différentes stratégies d'insertion (de type STDM). Les premiers résultats obtenus pour le filtre blanchissant sont très encourageants ; ils sont par contre très mitigés pour la STDM.

Références

- [1] RNRT, *ARTUS*, http://www.telecom.gouv.fr/rnrt/rnrt/projets/res_01_37.htm.
- [2] U. Zölzer, *DAFX - Digital audio effects*, Wiley, 2002.
- [3] D. Kirovski et H. Malvar, "Spread-spectrum watermarking of audio signals," *IEEE Transactions on Signal Processing*, vol. 51, pp. 1020–1033, Avril 2003.
- [4] A. LoboGuerrero, F. Marques, J. Lienard, et P. Bas, "Enhanced audio data hiding synchronization using non linear filters," *IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, vol. 2, pp. 885–888, mai 2004.
- [5] P. Bas, *Méthodes de tatouage d'images fondées sur le contenu*, Mémoire de Thèse, 2000.
- [6] S. Larbi, M. Jaidane, et N. Moreau, "A new wiener filtering based detection scheme for time domain perceptual audio watermarking," *IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, vol. 5, pp. 949–952, mai 2004.
- [7] J. Proakis, *Digital communications*, McGraw-Hill, 2001, 4ème édition.