

# Proposition d'une stratégie de fusion de données à trois niveaux pour la vérification d'identité

Patrick Verlinde\*    Dominique Genoud†    Guillaume Gravier    Gérard Chollet

Ecole Nationale Supérieure des Télécommunications / Département Signal  
46, rue Barrault  
75634 Paris Cedex 13, FRANCE

## ABSTRACT

The aim of this paper is to propose a strategy that uses data fusion at three different levels to gradually improve the performance of an identity verification system. In a first step temporal data fusion can be used to combine multiple instances of a single (mono-modal) expert to reduce its measurement variance. If the system performance is too low, one can try to improve it by fusing in a second step the results of multiple experts working on the *same* biometric modality. This supposes that their respective classification errors are (at least partially) uncorrelated. Finally, if the performance is still not good enough, one will be forced to pass on to the third step in which performance can be increased when using multiple experts working on different (biometric) modalities. These experts have to be chosen in such a way that the effect of adding the extra modalities is to increase the separation between the distributions of the different populations that have to be classified by the system. This kind of strategy allows a user to gradually tune the performance of an identity verification system to its needs while controlling the increase of investment costs<sup>1</sup>.

## 1. INTRODUCTION

L'objectif d'un système de vérification d'identité est de décider si une personne qui proclame une identité est bien la personne qu'elle prétend être. Chaque méthode implémentée, appelée ici expert, effectue un test d'hypothèse binaire. Ces tests génèrent deux types d'erreurs: un faux rejet (FR: rejeter un accès client) ou une fausse acceptation (FA: accepter un imposteur). Les performances d'un système de vérification d'identité sont en général caractérisées par les taux d'erreurs globaux mesurés pendant les tests: le taux de faux rejets [False Rejection Rate: FRR = (nombre de FR) / (nombre d'accès client)] et le taux de fausses acceptations [False Acceptance Rate: FAR = (nombre de FA) / (nombre d'accès imposteur)]. Une mesure unique peut être obtenue en combinant ces deux erreurs pour obtenir le taux d'erreurs global [Total Error Rate: TER = (nombre de FA + nombre de FR) / (nombre total d'accès)] ou par son complément, le

taux de réussite global [Total Success Rate: TSR = 1 - TER] [BC97]. Une vérification d'identité parfaite, FAR = 0 et FRR = 0, peut difficilement être obtenue en pratique. Cependant, une étude détaillée des tests d'hypothèse binaire [Tre68] montre que si l'on fait tendre FRR vers zéro en modifiant le seuil de décision, on augmente alors le taux de fausses acceptations, et inversement.

Les erreurs sont causées par l'imperfection des différents experts. On a donc des erreurs de mesure. Une solution, détaillée dans la section 2, pour remédier à ce problème consiste à combiner les résultats obtenus par un seul expert sur des instances successives sur plusieurs tests. On parlera alors de fusion temporelle. De cette manière on peut effectivement diminuer la variance des résultats, ce qui permet d'augmenter les performances du système de base sans avoir à modifier l'approche choisie. Si ces performances ne satisfont pas les besoins de l'utilisateur, il peut passer à l'amélioration suivante.

Le fait de disposer d'experts parfaits et de connaître les distributions statistiques des résultats obtenus par chaque méthode de vérification pour les deux classes (clients, imposteurs), permet d'utiliser la règle de décision bayésienne pour obtenir un fonctionnement du système au taux d'erreur minimal [DH73]. L'erreur de classification ne sera égale à zéro que dans le cas où les deux distributions ne se recouvrent pas du tout. Si cette condition n'est pas satisfaite, ce qui est le cas en général, on ne peut faire baisser les deux types d'erreurs (FAR, FRR) *simultanément* qu'en augmentant la taille de l'échantillon statistique, c'est-à-dire en utilisant plus de données d'entraînement, permettant ainsi de diminuer les variances des estimateurs des paramètres réels de la distribution. Très souvent on ne dispose malheureusement que d'un nombre limité de données d'entraînement, ce qui exclut cette possibilité théorique. Une autre possibilité pour faire diminuer ces deux taux d'erreurs simultanément consiste alors à fusionner les résultats de différents experts basés sur la même modalité. Ce type de fusion permet d'augmenter les performances si les erreurs commises par les experts ne sont pas totalement corrélées. L'amélioration des performances sera même d'autant plus grande que les résultats (et les erreurs) obtenus par les différents experts seront plus d'écourelés, l'apport d'information augmentant avec la d'écourelation. Dans la section 3, on montrera le cas d'un système mono-modal multi-expert basé sur la modalité biométrique vo-

\* Ecole Royale Militaire / Département ELTE, Bruxelles, Belgique, e-mail: patrick.verlinde@tele.rma.ac.be

† Institut Dalle Molle d'Intelligence Artificielle Perceptive, Martigny, Suisse, e-mail: genoud@idiap.ch

<sup>1</sup> Some of the results given in this paper were obtained in the framework of the CAVE-PICASSO and M2VTS European projects.

cale. Dans le cas spécifique d'une modalité vocale, on a ajouté la possibilité de différer la décision après fusion en cas de doute. Cette technique permet de diminuer les erreurs au prix d'une augmentation du nombre de non-décisions et d'un traitement supplémentaire. Ce mécanisme peut également être utilisé pour les autres systèmes de vérification d'identité, mais cette augmentation des performances induit des coûts supplémentaires. Si les performances du système après fusion mono-modale ne satisfont toujours pas l'utilisateur, il est possible de procéder à l'étape suivante.

Dans un troisième temps on peut espérer augmenter encore les performances en trouvant des caractéristiques qui améliorent la séparabilité des hypothèses testées, en augmentant la dimension de l'espace. Il est possible d'utiliser à cette fin des modalités biométriques supplémentaires. La discrimination entre les deux distributions sera d'autant plus facile que la corrélation entre les différentes modalités sera plus faible. On étudiera ce cas dans la section 4 pour des modalités biométriques vocales et visuelles.

## 2. FUSION MONO-MODALE MONO-MÉTHODE

### 2.1. Méthodes utilisées

Le premier pas dans la stratégie proposée, est d'utiliser une intégration temporelle afin d'augmenter les performances d'un système de vérification d'identité en diminuant la variance des mesures. Un exemple de cette technique est présentée dans [KMJS97]. Les auteurs de cet article partent, d'une part, du principe que les décisions d'accepter ( $\omega_1$ ) ou de rejeter ( $\omega_2$ ) la personne, sont basées sur la probabilité de classes *a posteriori*  $P(\omega_j|x_i)$ ,  $j = 1, 2$  et d'autre part, du fait que l'on dispose de plusieurs instances  $x_i$  de données biométriques. Dans ce cas, on peut se baser sur la fusion de l'information additionnelle. On suppose que les mesures  $x_i$  ont été acquises dans les mêmes conditions et qu'elles peuvent être considérées comme des mesures multiples qui ne diffèrent entre elles que par leur composante de bruit. Aussi les probabilités de classes *a posteriori* peuvent être considérées comme des estimations bruitées de la valeur nominale de cette probabilité, i.e.:

$$P(\omega_j|x_i) = P(\omega_j|x) + \epsilon(\omega_j|x_i).$$

Ainsi on peut obtenir de meilleures estimations  $\hat{P}(\omega_j|x)$  de cette probabilité *a posteriori* en combinant les estimations bruitées de manière linéaire ou en utilisant des statistiques de rang. Les différentes stratégies de fusion sont comparées entre elles en analysant leur probabilité d'erreur de manière théorique. Les méthodes de fusion étudiées sont les suivantes:

#### Moyenne arithmétique

$$\hat{P}(\omega_j|x) = \frac{1}{R} \sum_{i=1}^R P(\omega_j|x_i).$$

#### Règle maximale

$$\hat{P}(\omega_j|x) = \max_{i=1}^R P(\omega_j|x_i).$$

#### Règle minimale

$$\hat{P}(\omega_j|x) = \min_{i=1}^R P(\omega_j|x_i).$$

#### Règle médiane

$$\hat{P}(\omega_j|x) = \text{med}_{i=1}^R P(\omega_j|x_i).$$

## 2.2. Résultats

On démontre que la réduction de la variance est la plus spectaculaire pour la moyenne arithmétique près de la frontière de décision. Les réductions de la variance pour les statistiques de rang dépendent, entre autres, du nombre d'instances  $R$  et de l'ordre du rang et, bien qu'elles soient moins efficaces que dans le cas de la moyenne arithmétique, elles sont plus robustes par rapport aux valeurs extrêmes des résultats de mesure. La table 1 montre une réduction des taux d'erreurs jusqu'à 40%. Le gain en performance est monotone au début d'une séquence, mais tend à saturer après intégration des premières observations. Ces résultats ont été obtenus en utilisant un expert de vérification visuelle, basé sur une corrélation robuste optimisée. Pour ces expériences, l'expert utilise uniquement la modalité image de la base de données multi-modale M2VTS [PV97] qui contient 37 personnes.

**Table 1:** EER en fonction du nombre d'instances

Expert	1	2	3	4	5	6
Moyenne	6.9	4.6	4.5	3.7	4.0	4.0
Médian	6.9	4.8	4.2	4.3	4.2	4.0

## 3. FUSION MONO-MODALE MULTI-MÉTHODE

### 3.1. Méthodes utilisées

La combinaison de plusieurs experts sur une même modalité biométrique est très utilisée. On peut citer par exemple [Ben92, XKS92].

A titre d'exemple plus élaborée, les expériences décrites dans [GBGC96], basées sur plusieurs méthodes de vérification du locuteur en mode dépendant du texte, sont décrites. Les trois experts utilisés acceptent en entrée un ensemble de vecteurs de 12 coefficients cepstraux LPC extraits d'un signal de parole. Ces experts sont:

1. *Dynamic Time Warping (DTW)* L'algorithme de DTW consiste à effectuer une comparaison dynamique entre une matrice de référence et une matrice de test. Le résultat est une mesure de distance entre le test et la référence.
2. *Statistique du second ordre (SSO)* Une matrice de covariance  $X = \frac{1}{M} \sum_{t=1}^{t=M} X_t X_t^T$  de la séquence de parole de référence est créée. Une matrice de covariance  $Y$  de la séquence de test est également créée. Une mesure de sphéricité symétrique  $\text{Sym}(X, Y)$  est ensuite effectuée.

3. *Modèles de Markov cachés (HMM)* Deux sortes de modèles HMM sont créées pour chaque mot à tester: (1) **Un modèle du monde**, créée à partir d'une base de données (Polyphone [CCCL95]). Les modèles de monde sont identiques pour tous les clients. (2) **Un modèle du client**, qui utilise comme paramètres initiaux le modèle du monde, et dont on réestime les paramètres pour chaque locuteur avec les données de celui-ci. La mesure utilisée consiste à calculer le log du rapport de vraisemblance locuteur/monde pour un segment de parole donnée.

**Détermination des seuils** La décision d'acceptation ou de rejet d'un locuteur s'effectue par comparaison du score de chaque méthode à un *seuil a priori*. Le seuil est déterminé par la méthode de Furui. Ce seuil est robuste lorsque très peu de données d'entraînement sont à disposition [Fur81].

**Combinaison des décisions** Afin d'améliorer la décision globale de notre système, nous avons combiné les décisions de chaque méthode (DTW, SSO, HMM). La combinaison utilisée ici est un système majoritaire pondéré avec seuil de doute. La décision par rapport au seuil *a priori* est pondérée (entre 0 et 1) en fonction de la distance seuil/score, ce qui revient à établir une confiance dans cette décision. Le système compare la moyenne des confiances des méthodes qui ont obtenu la majorité à un seuil (seuil de doute). Si cette moyenne est supérieure au seuil, on prend la décision de la majorité, sinon il y a doute.

### 3.2. Résultats

La base de données Polycode [GC95] utilisée pour ces expériences est constituée de 25 locuteurs. La table 2 montre que la fusion des décisions des trois experts est meilleure que le meilleur des experts. Pris dans un cadre plus formel, les trois experts utilisés ne sont pas fondamentalement différents, la DTW et les SSO étant des cas "dégénérés" de la modélisation HMM. Cependant la modélisation des locuteurs effectuée par chaque expert n'étant pas parfaite, ce sont parfois les algorithmes nécessitant moins de paramètres à estimer qui s'avèrent les plus robustes. Parmi les facteurs qui entachent l'estimation correcte d'un modèle, le manque de données d'entraînement et la variabilité intra-locuteur sont certainement les plus importantes.

**Table 2:** Performance des experts et de la fusion

Expert	FRR% (200 tests)	FAR% (1800 tests)
DTW	23.5	7.67
SOSM	14.0	5.28
HMM L/R	5.53	2.72
Fusion	2.0	2.72

Malgré l'amélioration des performances due à la fusion, les erreurs résiduelles peuvent, dans certaines applications, s'avérer inacceptables, et il est possible, dans ce cas,

d'éviter de prendre de mauvaises décisions en rejetant les cas douteux. La table 3 nous indique le comportement des performances du système lors du changement de seuil de doute. Ce seuil permet de régler ainsi les taux de faux rejets et fausses acceptations selon le niveau de sécurité requis par l'application. Les cas douteux peuvent être résolus en demandant des informations supplémentaires au locuteur.

**Table 3:** Performance en changeant le seuil de doute

Seuil de doute	FRR% (200 tests)	FAR% (1800 tests)	Doute% (2000 tests)
0.2	2.0	2.72	0.0
0.5	2.0	2.33	0.83
0.7	2.0	1.11	10.11
0.8	1.0	0.89	20.2

## 4. FUSION MULTI-MODALE MULTI-MÉTHODE

### 4.1. Méthodes utilisées

Le troisième pas dans la stratégie de développement proposée est la fusion de multiples experts utilisant plusieurs modalités (biométriques). Cette approche a, entre autre, été utilisée dans les références suivantes [BF95, DBB, JLGW97].

Pour montrer les possibilités de la fusion multi-modale, nous avons testé différents modules de fusion et nous avons comparé les résultats à ceux obtenus avec les experts individuels, sur la base de données multi-modale M2VTS [PV97] avec cinq experts:

**UCL Pro** expert visuel basé sur les profils (Université Catholique de Louvain);

**UCL Grey** expert visuel basé sur les profils et les niveaux de gris (Université Catholique de Louvain);

**RMA Pro** expert visuel basé sur les profils (Ecole Royale Militaire);

**UOS Cor** expert visuel frontal basé sur une méthode de corrélation robuste (University of Surrey);

**RMA Ar** expert vocal basé sur une méthode autorégressive (Ecole Royale Militaire).

On a testé les modules de fusion suivants: méthodes par vote (ET, OU, Majoritaire) [Das94], classificateurs linéaire et quadratique [DH73] et le perceptron multicouches (MLP) [Bis95].

### 4.2. Résultats

Les résultats expérimentaux se trouvent dans la table 4. Ces résultats, exprimés en taux d'égales erreurs, ont été obtenus en utilisant les seuils *a posteriori* et représentent donc les meilleures performances possibles.

Après ce premier test, on a comparé la meilleure méthode avec un MLP (caractérisé par 25 neurones logsig dans la

**Table 4:** Taux d'égales d'erreurs avant et après fusion

Expert	EER%
UCL Pro	08.0
UCL Grey	07.0
RMA Pro	17.0
UOS Cor	12.0
RMA Ar	02.0
Méthode par vote ET	12.5
Méthode par vote OU	06.5
Méthode par vote MAJ	05.0
Classificateur linéaire	02.5
Classificateur quadratique	01.0

couche cachée et un neurone logsig dans la couche de sortie), sur des données de test en utilisant les seuils calculés *a posteriori* (selon le critère EER). Les résultats se trouvent dans la table 5. On constate que la fusion de données multi-modale permet effectivement d'augmenter les performances au delà des performances du meilleur expert.

**Table 5:** Comparaison des erreurs de vérification

Méthode	FRR%	FAR%
Classificateur quadratique	0.0	4.1
MLP	0.0	0.4

## 5. CONCLUSIONS

La fusion de données peut être utilisée au sein d'un système de vérification d'identité à différents niveaux. Le premier niveau est la fusion temporelle des résultats obtenus par un seul expert, et donc sur une seule modalité, afin de réduire la variance de mesure. Un deuxième niveau est atteint en fusionnant les résultats obtenus par différents experts sur la même modalité, ceci afin de minimiser les erreurs de classification en jouant sur la décorrélation des erreurs des différents experts. Le troisième niveau d'application est de diminuer encore les erreurs de classification en essayant d'accroître la distance entre les distributions des deux populations en augmentant la dimension de l'espace par l'utilisation de plusieurs modalités (biométriques) les plus décorréées possible. Nos travaux actuels sont axés vers une combinaison de ces trois utilisations de la fusion de données dans une seule stratégie séquentielle, qui a le grand avantage d'augmenter les performances d'un système de vérification d'identité de manière graduelle.

## BIBLIOGRAPHIE

- [BC97] F. Bimbot and G. Chollet. Assessment of speaker verification systems. In *Handbook of Standards and Resources for Spoken Language Systems*. Mouton de Gruyter, 1997.
- [Ben92] Y. Bennani. Text-independent talker identification system combining connectionist and conventional models. In S. Y. Kung et al., editor, *Neural Networks for Signal Processing, Vol.2*. IEEE Service Center Press, 1992.

- [BF95] R. Brunelli and D. Falavigna. Person identification using multiple cues. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(10):955–966, October 1995.
- [Bis95] C. M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, Oxford UK, 1995.
- [CCCL95] G. Chollet, J.L. Cochard, A. Constantinescu, and Ph. Langlais. Swiss french polyphone and polyvar: telephone speech databases to study intra and inter speaker variability. Technical report, IDIAP, 1995.
- [Das94] B. V. Dasarathy. *Decision Fusion*. IEEE Computer Society Press, 1994.
- [DBB<sup>+</sup>97] B. Duc, E. Bigün, J. Bigün, G. Maître, and S. Fischer. Fusion of audio and video information for multi modal person authentication. *Pattern Recognition Letters*, 18(9):835–843, 1997.
- [DH73] R. O. Duda and P. E. Hart. *Pattern Classification and Scene Analysis*. John Wiley & Sons, New York, 1973.
- [Fur81] S. Furui. Cepstral analysis technique for automatic speaker verification. *IEEE Trans. on ASSP*, ASSP-29(2):254–272, 1981.
- [GBGC96] D. Genoud, F. Bimbot, G. Gravier, and G. Chollet. Combining methods to improve speaker verification decision. In *Proceedings of The Fourth International Conference on Spoken Language Processing*, 1996.
- [GC95] Dominique Genoud and Gérard Chollet. Polycode a verification database. Technical report, IDIAP, CH-1920 Martigny, 1995.
- [JLGW97] P. Jourlin, J. Luetttin, D. Genoud, and H. Wassner. Acoustic-labial speaker verification. *Pattern Recognition Letters*, 18(9):845–852, 1997.
- [KMJS97] J. Kittler, G. Matas, K. Jonsson, and M. U. Ramos Sánchez. Fusion of audio and video information for multi modal person authentication. *Pattern Recognition Letters*, 18(9):845–852, 1997.
- [PV97] S. Pigeon and L. Vandendorpe. The M2VTS multi-modal face database. In *Proceedings of the First International Conference on Audio- and Video-based Biometric Person Authentication (AVBPA'97)*, Lecture Notes in Computer Science. Springer Verlag, 1997.
- [Tre68] H. L. Van Trees. *Detection, Estimation and Modulation Theory*, volume 1. John Wiley & Sons, New York, 1968.
- [XKS92] L. Xu, A. Krzyzak, and C.Y. Suen. Methods of combining multiple classifiers and their applications to handwriting recognition. *IEEE Transactions on Systems, Man and Cybernetics*, 22(3):418–435, May/June 1992.