

ANALYSE EXPLORATOIRE & PREPROCESSING

		Valeur	Valeur en %
Vérification	Colonne		
Valeurs manquantes	designation	0	0.000000
	description	29800	35.093504
	productid	0	0.000000
	imageid	0	0.000000
	prdtypecode	0	0.000000
Doublons	Designation	2651	3.121909
	Description	7610	13.807243
	Productid	0	0.000000
	Imageid	0	0.000000
	Prdtypecode	84889	99.968204
Unicité	Designation	82265	96.878091
	Description	47506	55.944698
	Productids	84916	100.000000
	Imageids	84916	100.000000
	Prdtypecode	27	0.031796

Tableau : Résumé de la qualité des données : valeurs manquantes, doublons et unicité des colonnes clés

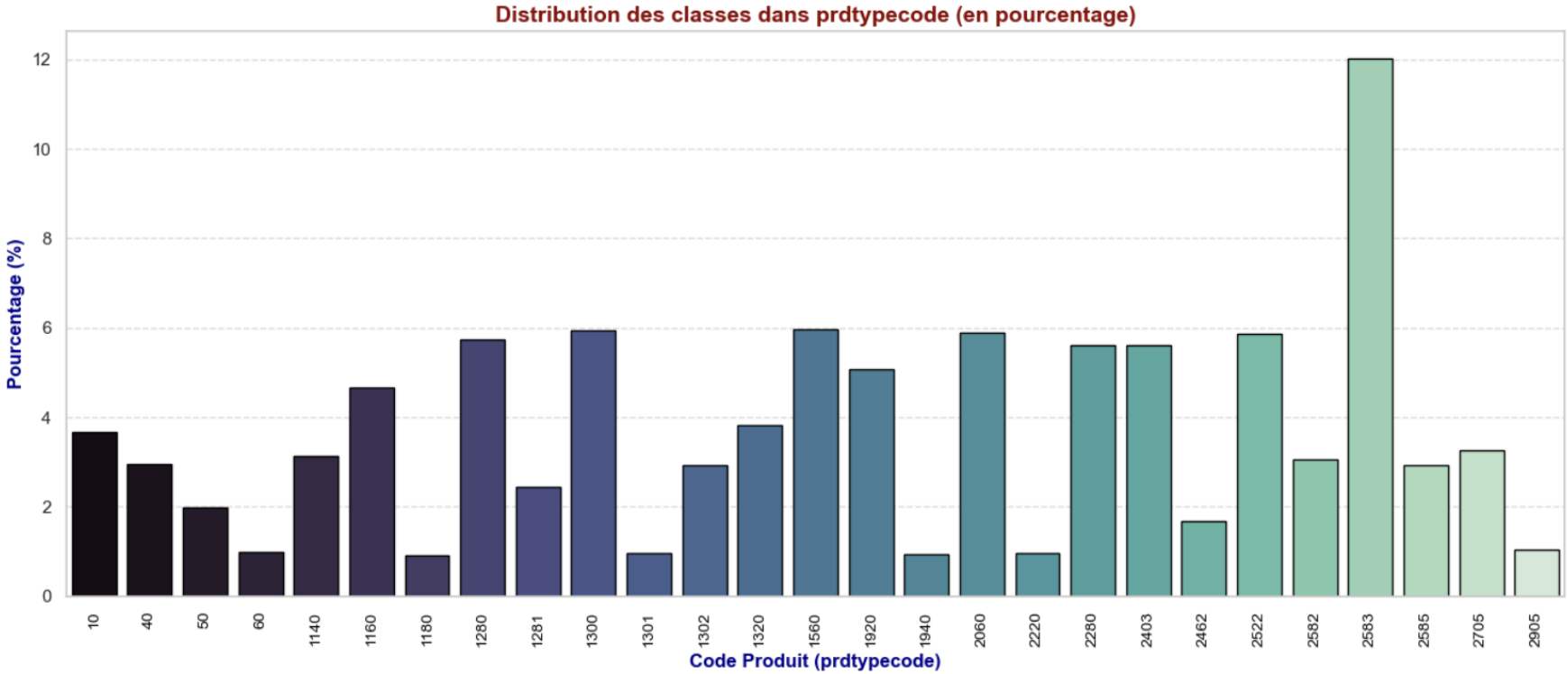


Figure 1 : Distribution catégories de produits (prdtypecode) dans le jeu de données Rakuten

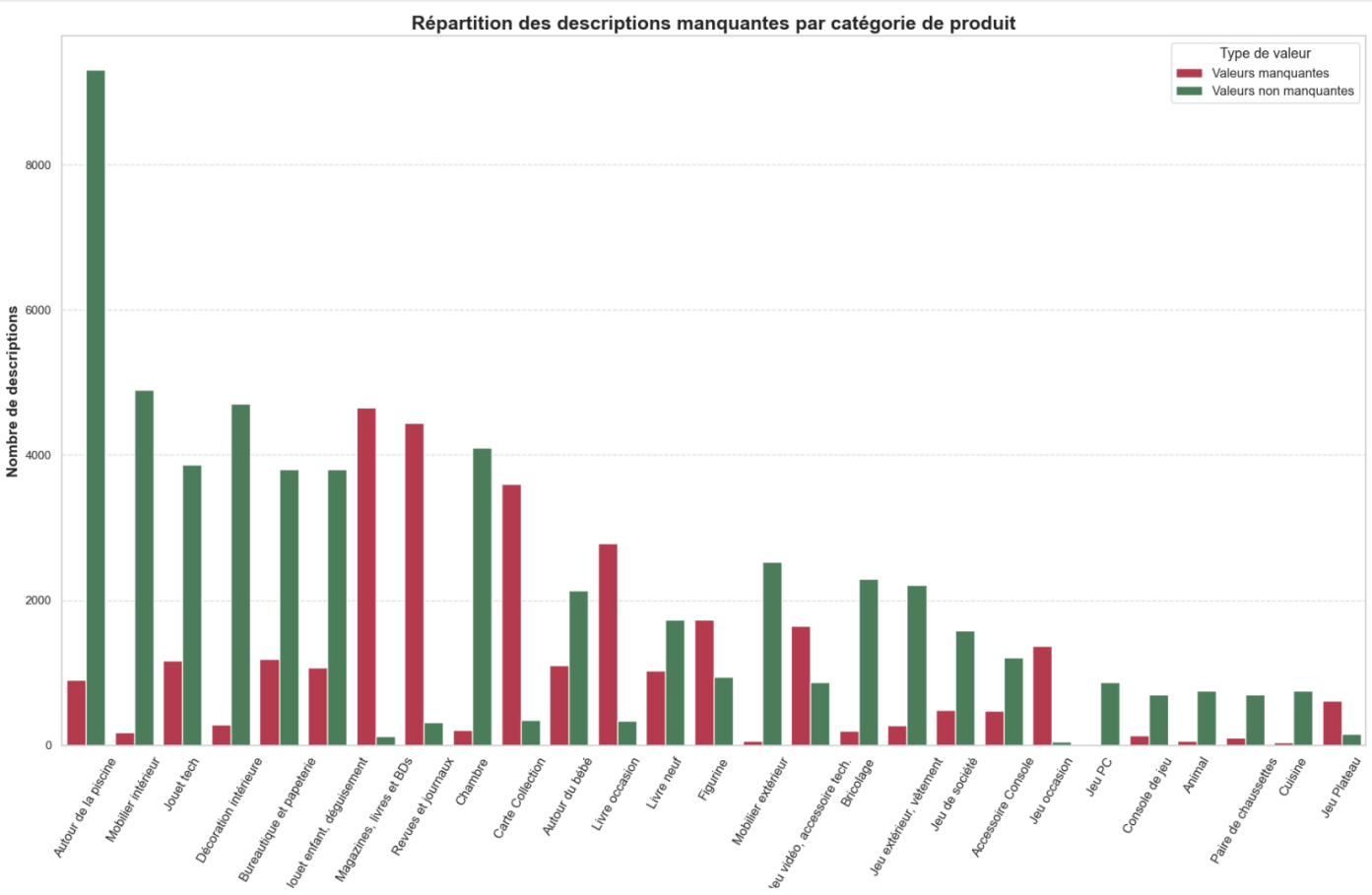


Figure 2 : Répartition des valeurs manquantes dans la colonne description selon la catégorie de produit

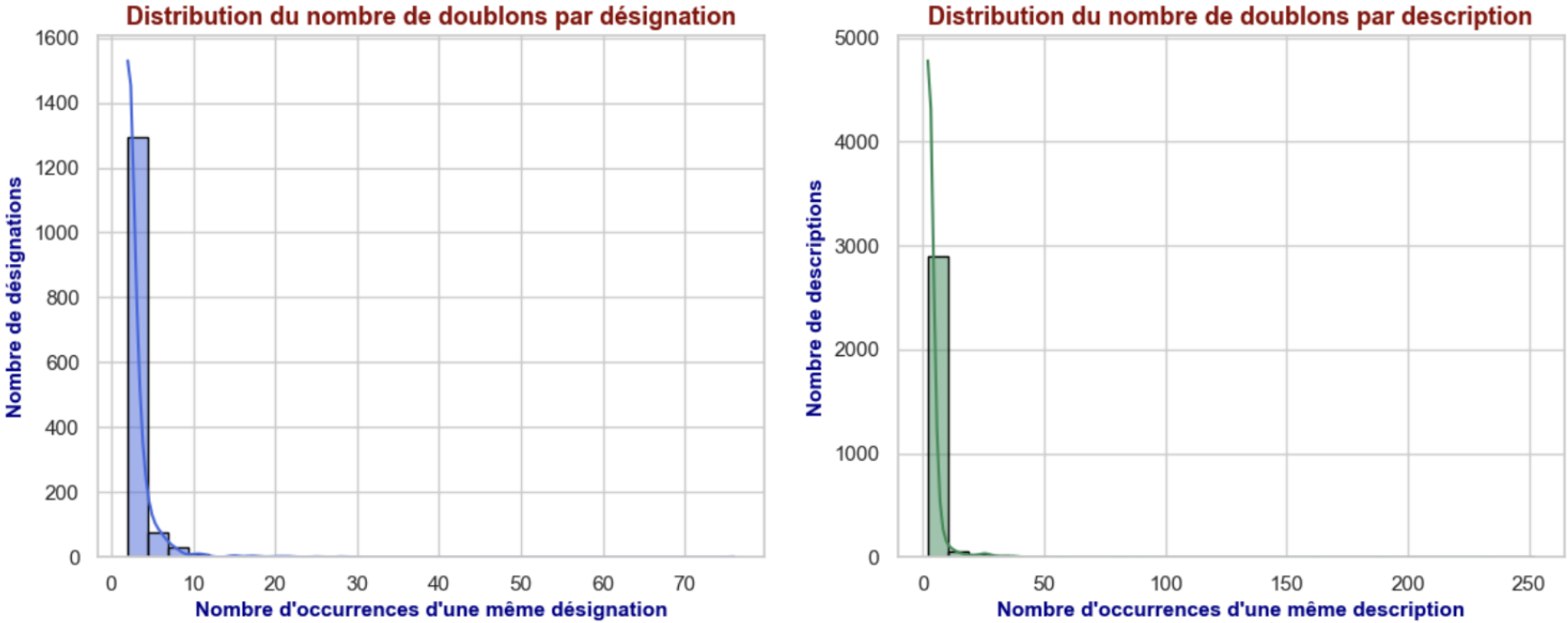


Figure 3 : Distribution du nombre de doublons dans les colonnes textuelles designation et description

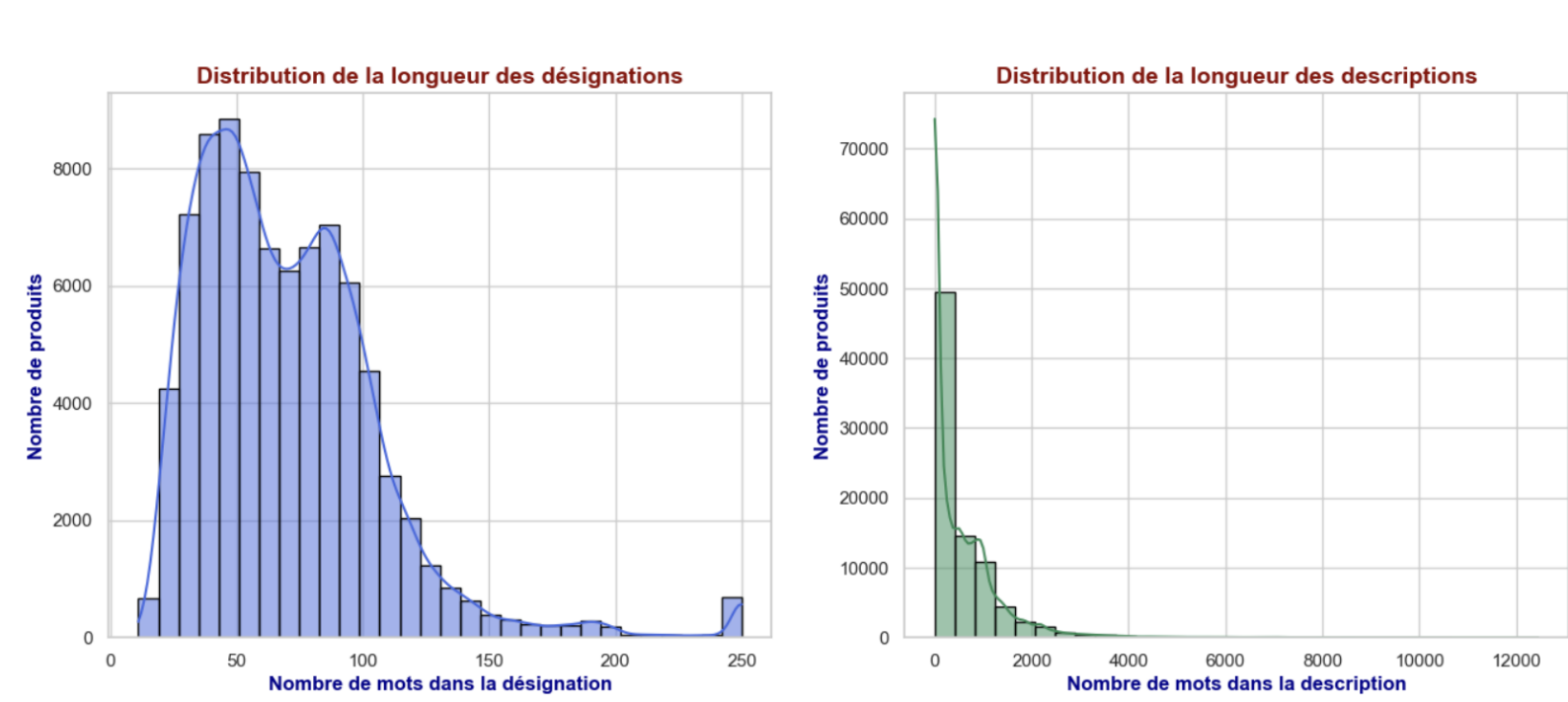




Figure 7: Nuage de mots combiné : contenu textuel des produits Rakuten **sans préprocessing**

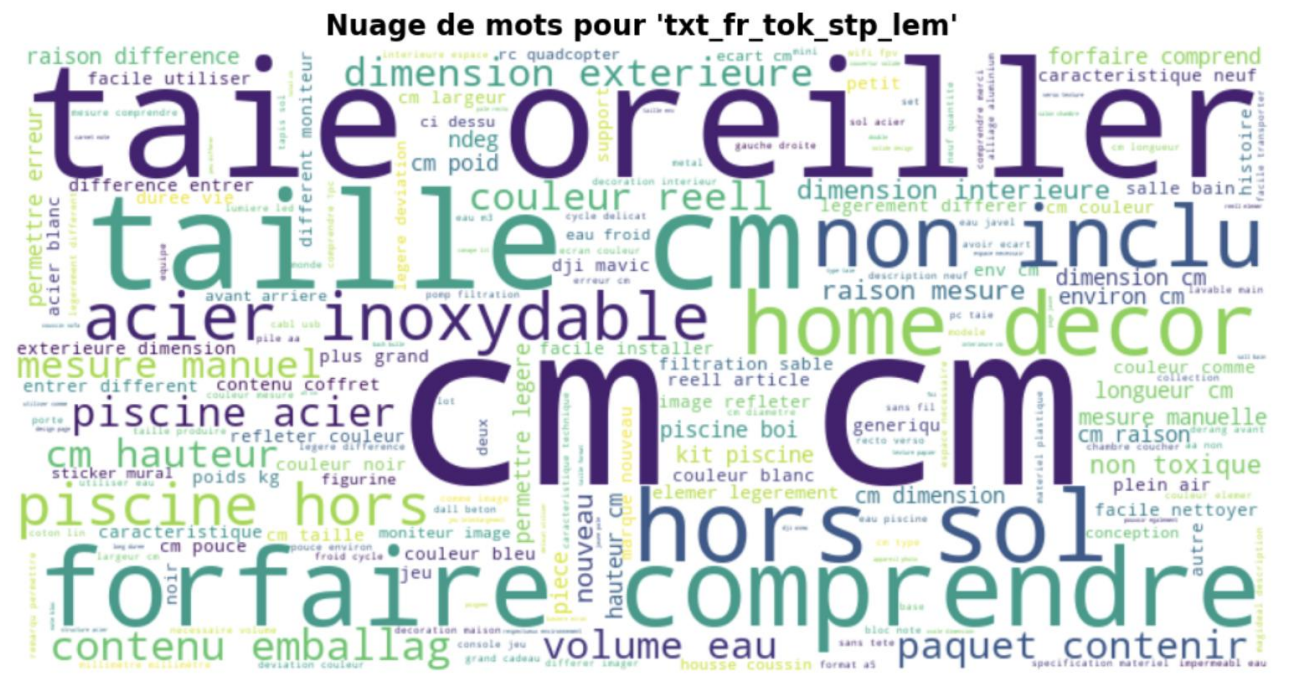


Figure 8: Nuage de mots globale et combiné des textes traduits en français **après preprocessing**

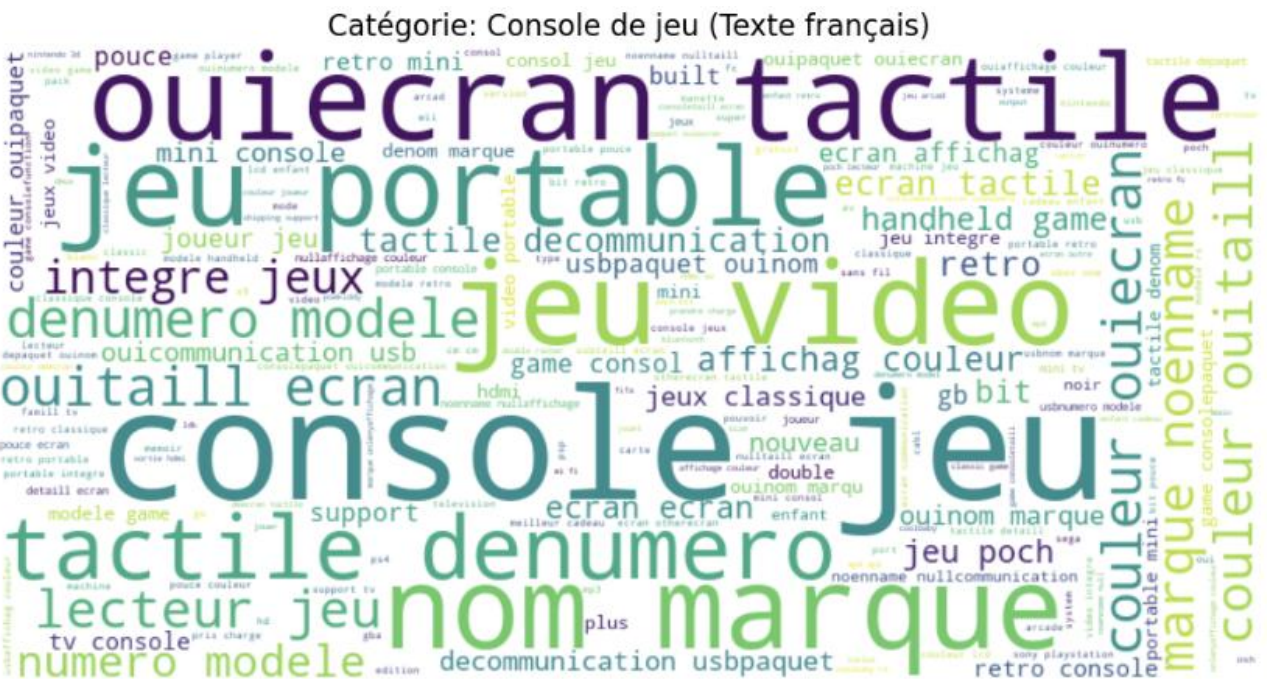


Figure 9: Nuage de mots globale et combiné après preprocessing la catégorie produit 'Console de jeu' : corpus multilingue vs français

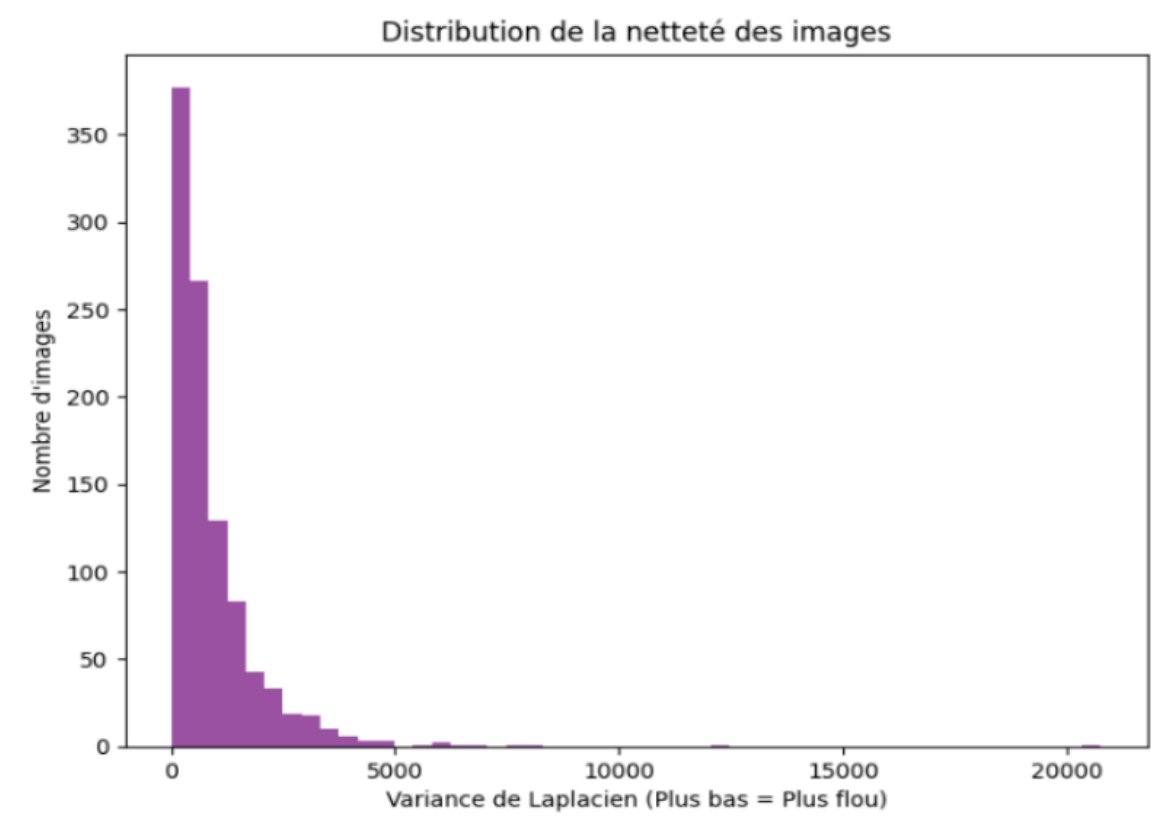


Figure 10: Distribution de la netteté des images mesurée par la variance du Laplacien

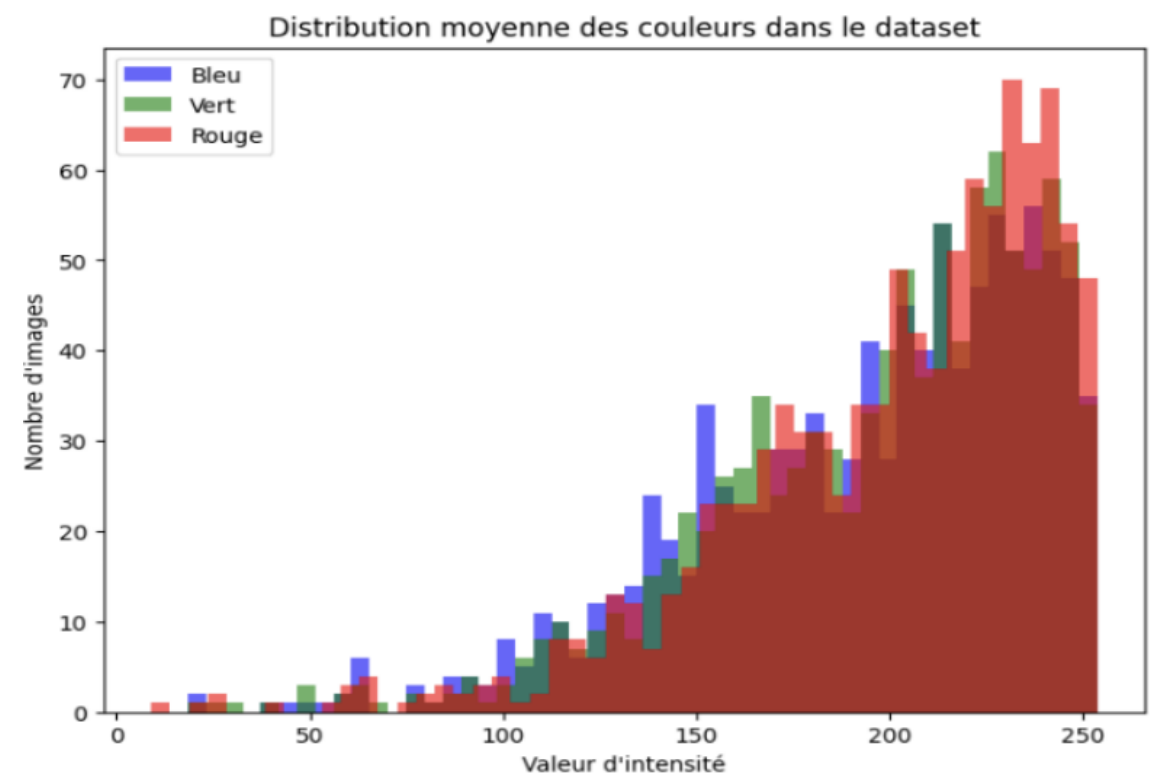
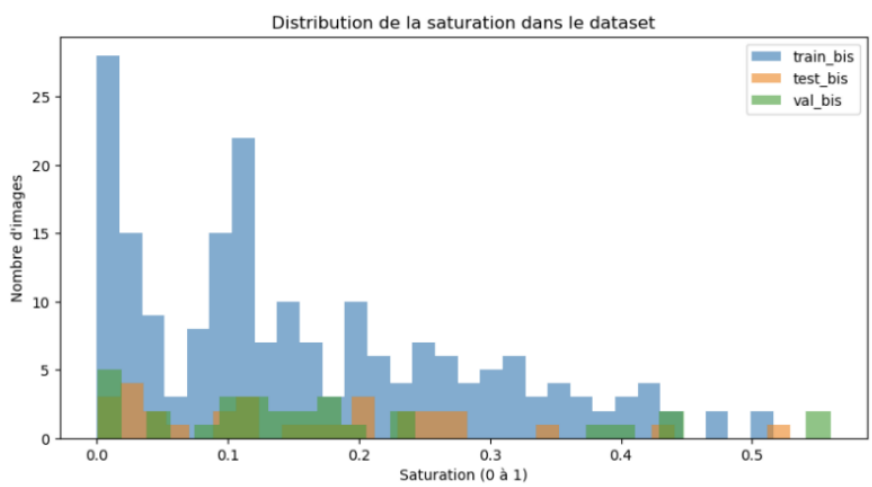
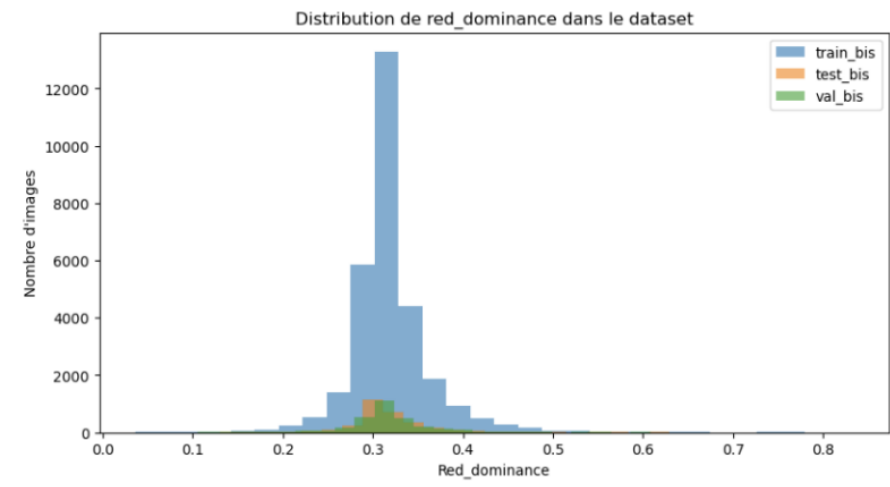


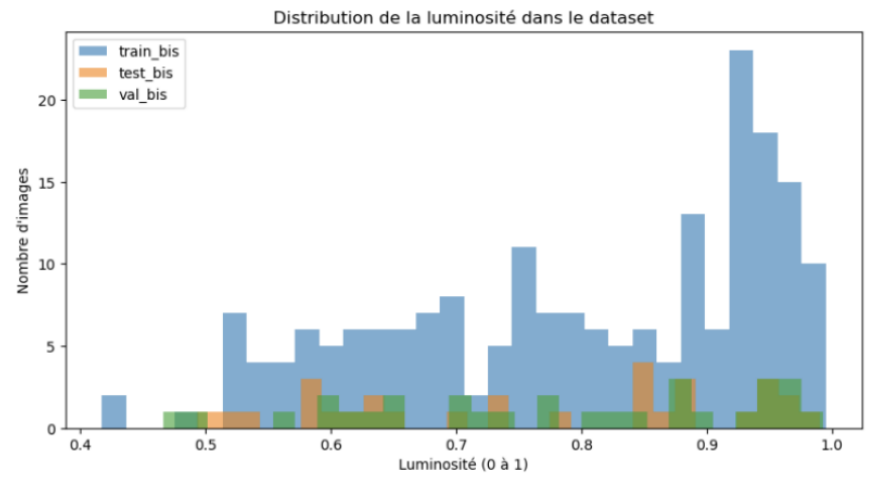
Figure 11: Distribution des intensités moyennes des canaux de couleur (RVB) dans les images du dataset



Dataset	Saturation moyenne
train_bis	0.201
test_bis	0.201
val_bis	0.224



Dataset	Red_dominance moyenne
train_bis	0.318
test_bis	0.320
val_bis	0.318



Dataset	Luminosité moyenne
train_bis	0.782
test_bis	0.766
val_bis	0.775

Figure 12: Analyse des propriétés colorimétriques des images : saturation (gauche), dominance du rouge (milieu) et luminosité (droite)

MODELISATION

Classification Report sur Validation:

	precision	recall	f1-score	support
Accessoire Console	0.74	0.83	0.78	166
Animal	0.82	0.87	0.84	82
Autour de la piscine	0.98	0.97	0.97	989
Autour du bébé	0.83	0.81	0.82	324
Bricolage	0.79	0.82	0.80	244
Bureautique et papeterie	0.93	0.93	0.93	496
Carte Collection	0.89	0.89	0.89	396
Chambre	0.91	0.92	0.91	427
Console de jeu	0.87	0.86	0.86	83
Cuisine	0.88	0.89	0.88	81
Décoration intérieure	0.82	0.77	0.80	497
Figurine	0.73	0.80	0.77	267
Jeu PC	0.92	0.99	0.96	87
Jeu Plateau	0.49	0.63	0.55	76
Jeu de société	0.60	0.56	0.57	207
Jeu extérieur, vêtement	0.85	0.79	0.82	249
Jeu occasion	0.72	0.74	0.73	141
Jeu vidéo, accessoire tech.	0.66	0.66	0.66	251
Jouet enfant, déguisement	0.76	0.66	0.71	487
Jouet tech	0.93	0.89	0.91	505
Livre neuf	0.67	0.72	0.69	276
Livre occasion	0.49	0.56	0.52	311
Magazines, livres et BDs	0.78	0.74	0.76	475
Mobilier extérieur	0.72	0.77	0.74	257
Mobilier intérieur	0.84	0.84	0.84	505
Paire de chaussettes	0.90	0.93	0.91	80
Revue et journaux	0.84	0.84	0.84	473
accuracy			0.82	8432
macro avg	0.79	0.80	0.80	8432
weighted avg	0.82	0.82	0.82	8432

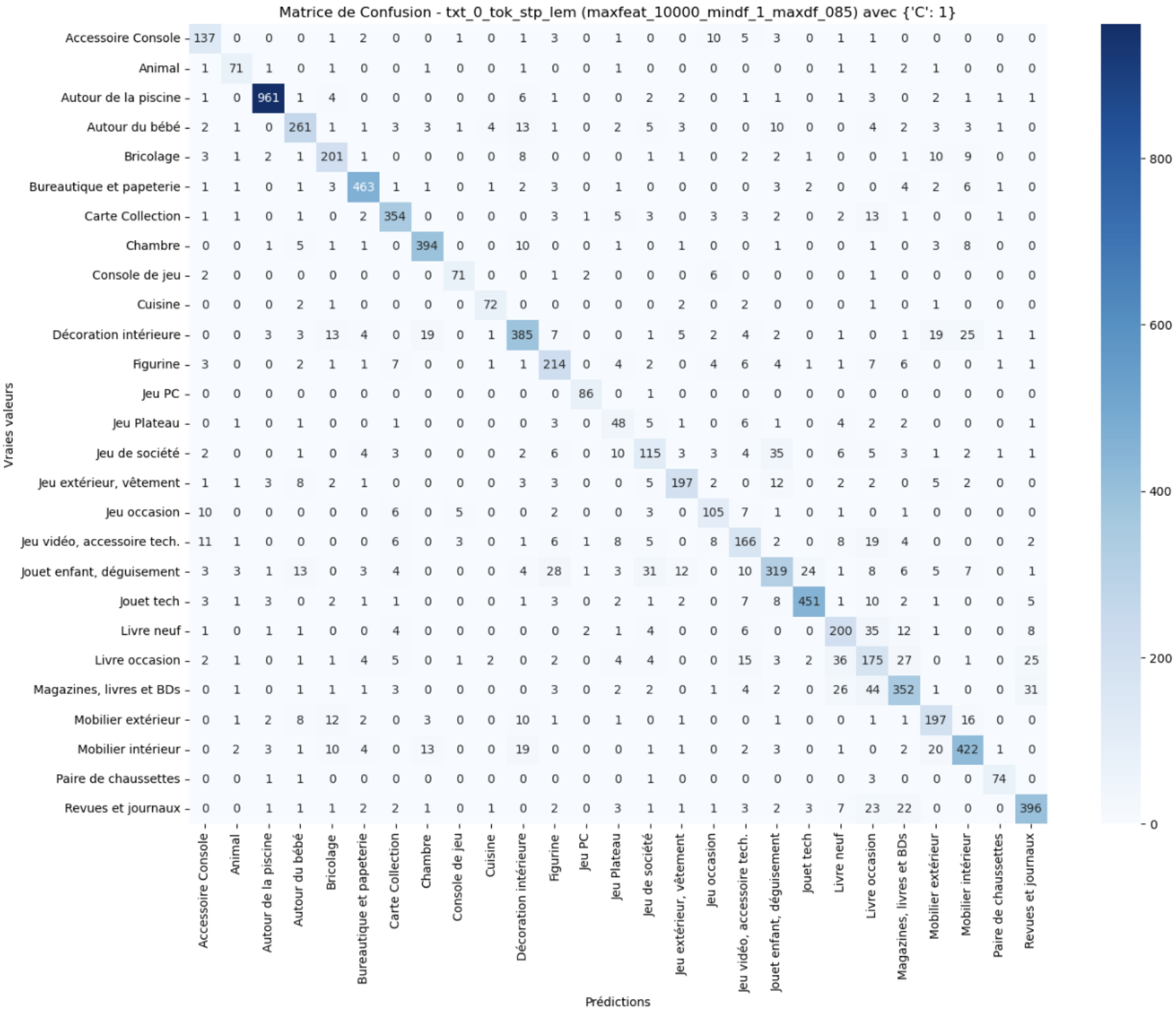


Figure 13: Performance du meilleur modèle Machine Learning : LinearSVC (TF-IDF, class_weight=balanced, C=1) sur le jeu de validation avec rapport de classification et matrice de confusion

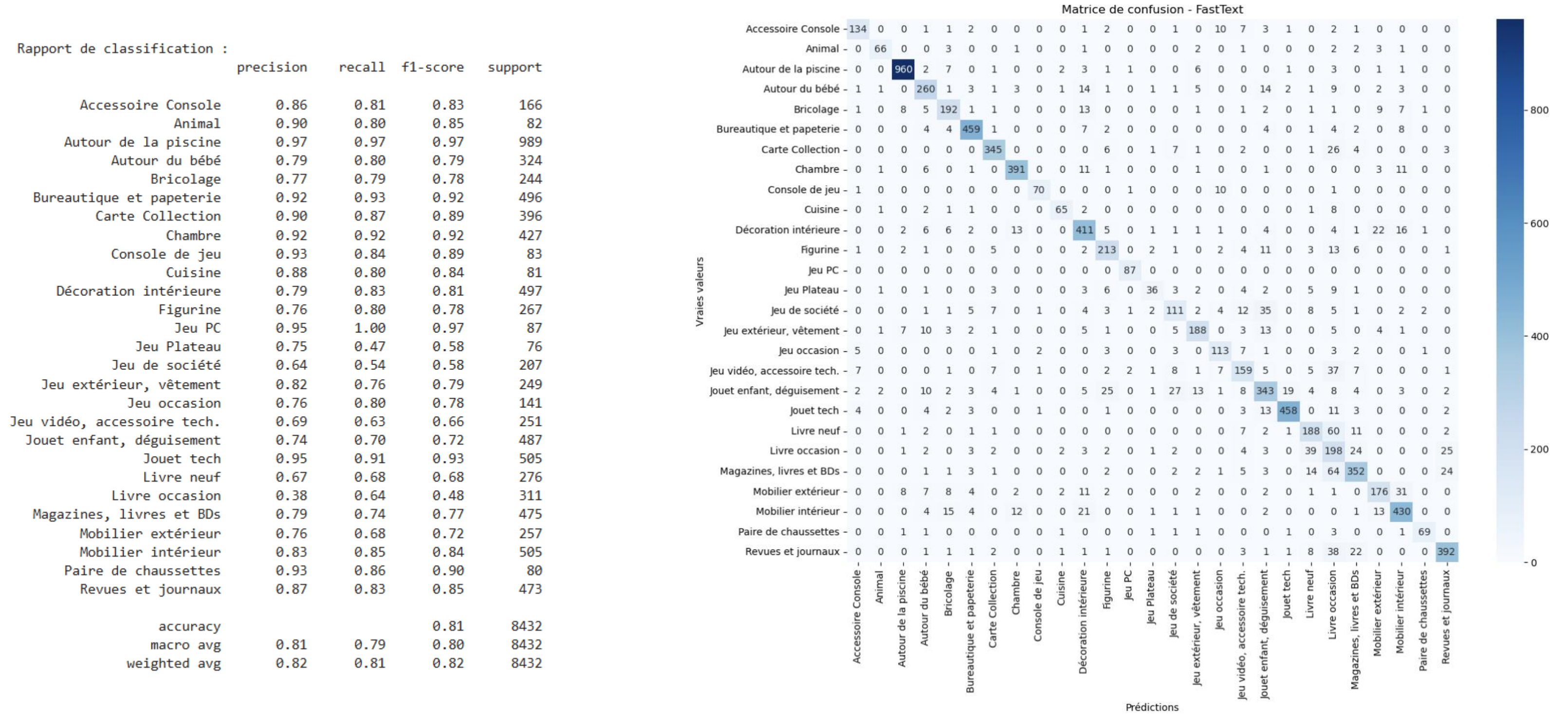


Figure 14: Performance du meilleur modèle FastText supervised (lr=0.5, dim=100, epoch=100) sur le jeu de validation avec rapport de classification et matrice de confusion

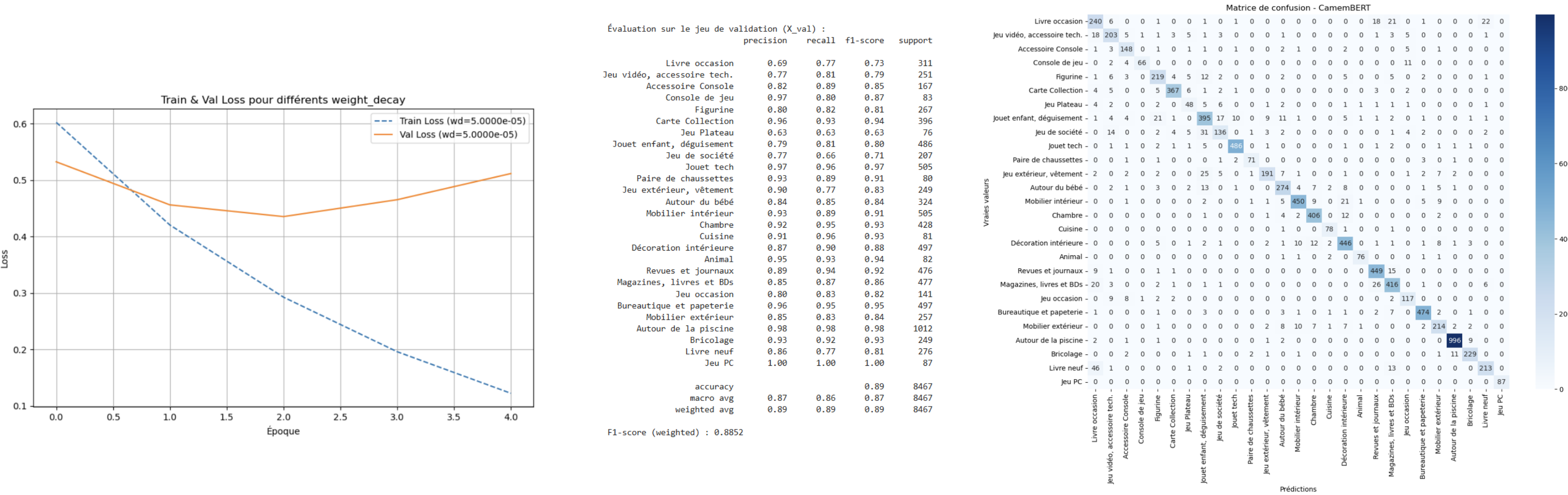


Figure 15: Performance du meilleur modèle CamemBERT (lr = 5e-5, weight_decay = 0.001, batch_size=16) sur le jeu de validation : perte (gauche), rapport de classification (milieu), matrice de confusion (droite)

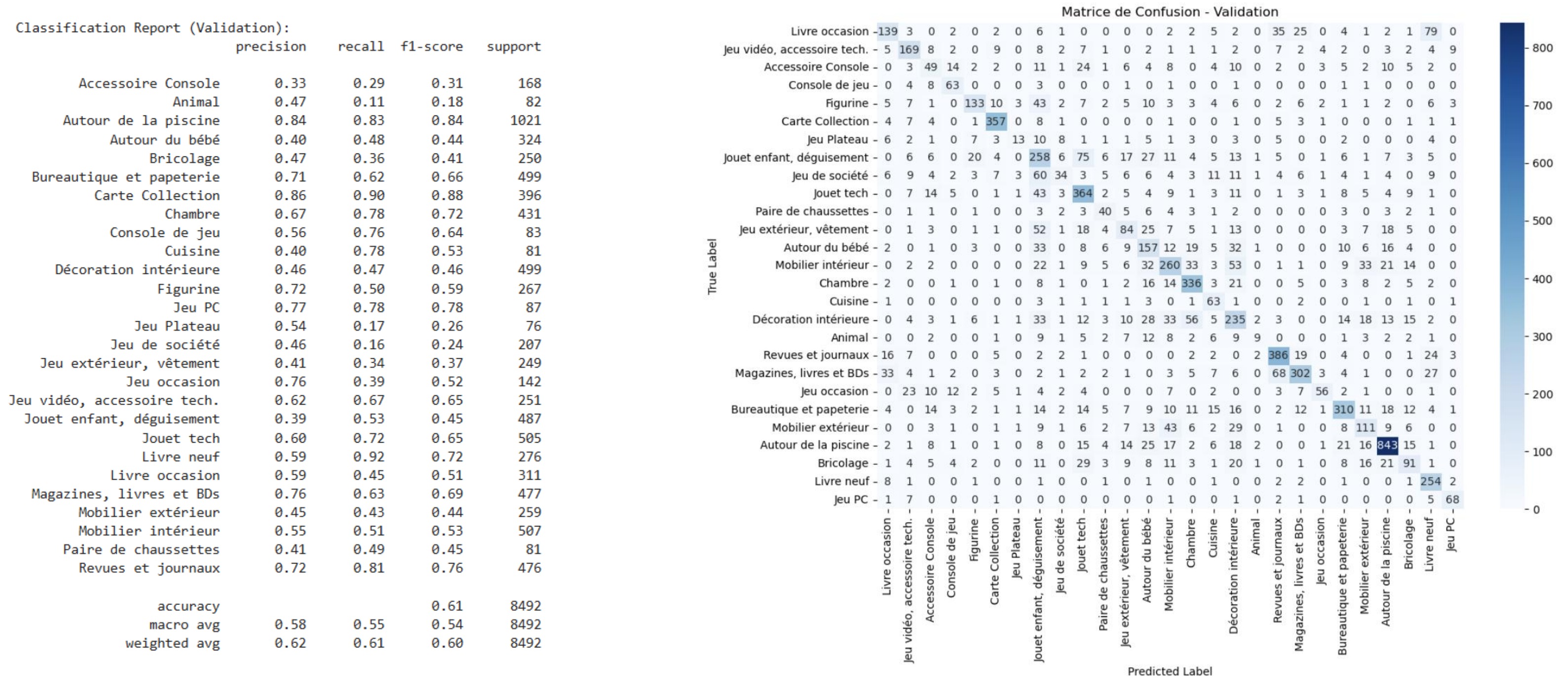


Figure 16: Performance du meilleur modèle RESNET50 (avec data-augmentation) sur le jeu de validation : perte (gauche), rapport de classification (milieu), matrice de confusion (droite)

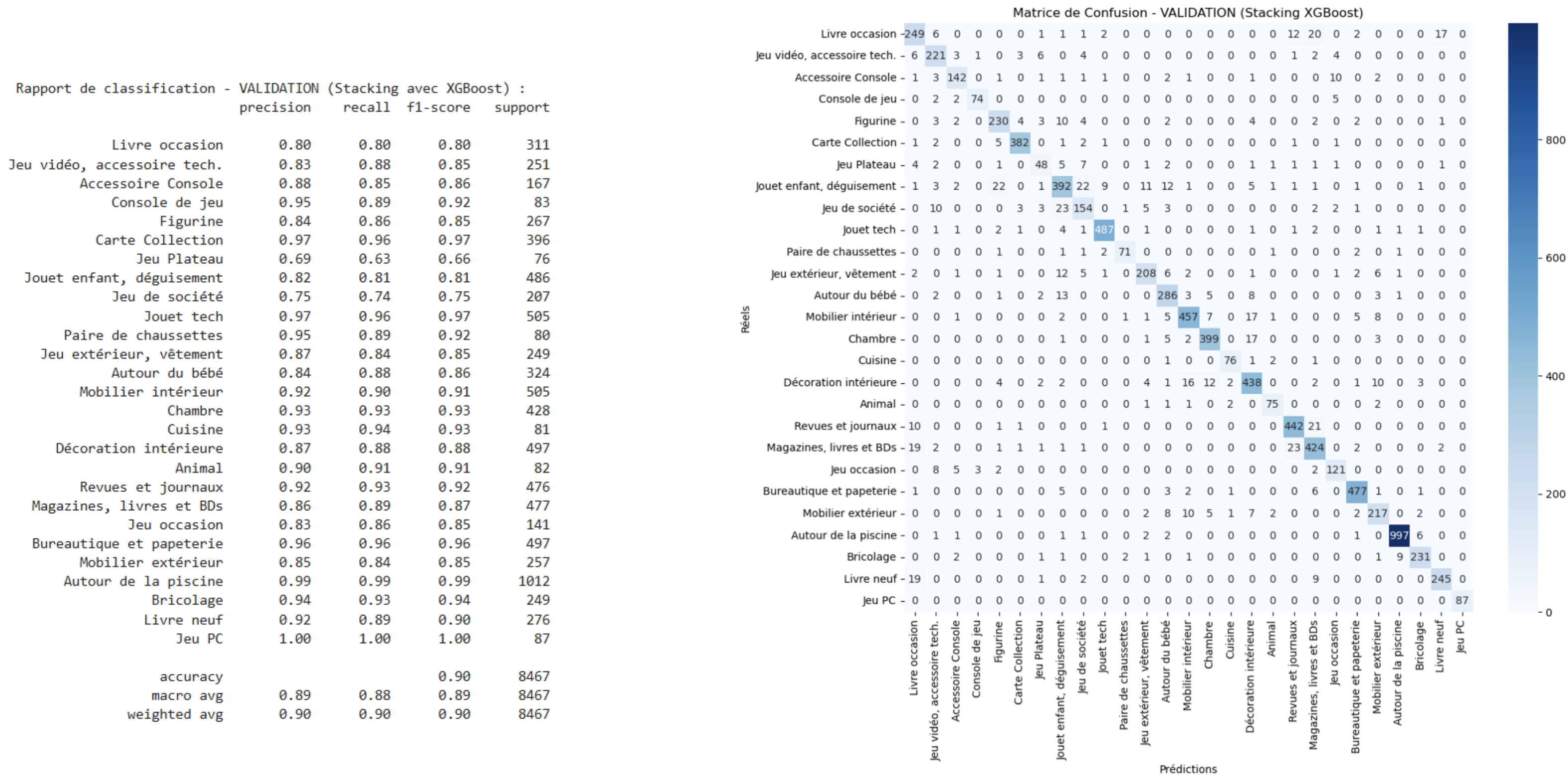


Figure 17: Performance du meilleur modèle multimodale (Stacking Classfier, Level-0: SVC sur probabilités CamemBERT, RandomForest sur probabilités RESNET50), Level-1: XGBoost) sur le jeu de validation : perte (gauche), rapport de classification (milieu), matrice de confusion (droite)