

# Proposition d'atelier pour PFIA2026

Contributeurs : Faten Chaieb-Chakchouk et Djallel DILMI

## Titre proposé :

Des modèles de langage aux modèles visuels : mécanismes d'attention (Vaswani et al. 2017) et applications en vision par ordinateur

## Description courte (2 phrases)

Ce tutoriel explore les mécanismes d'attention, depuis leur émergence dans les modèles de langage (LLM) jusqu'à leur adaptation en vision par ordinateur. Il met l'accent sur les fondements mathématiques, les intuitions conceptuelles et les applications pratiques en vision.

## Description longue (2 paragraphes)

Les mécanismes d'attention ont profondément transformé l'intelligence artificielle moderne, en particulier avec les modèles de langage de grande taille (LLM) basés sur les Transformers. Ce tutoriel présente une introduction progressive aux concepts fondamentaux de l'attention, non pas comme une simple opération calculatoire, mais comme un produit scalaire dynamique entre représentations, avec un lien direct avec des méthodes classiques telles que l'analyse canonique des corrélations (CCA). Cette approche permet de comprendre pourquoi l'attention capture efficacement les dépendances complexes et les relations contextuelles.

Dans une seconde partie, le tutoriel se concentre sur l'adaptation de ces mécanismes au domaine de la vision par ordinateur. Seront abordés les Vision Transformers -ViT- (Dosovitskiy et al. 2020), l'attention spatiale et canal -SE, CBAM, Mamba- (Wang et al. 2024, Gu et al. 2023), ainsi que les architectures hybrides CNN + attention. Une session pratique permettra aux participants de manipuler des modèles d'attention appliqués à des données visuelles, de visualiser les cartes d'attention et de comprendre leurs avantages et limites dans des tâches de classification et de détection.

## Déroulé détaillé (3h)

Étape / Durée	Contenu
Introduction (15–20 min)	<ul style="list-style-type: none"><li>Évolution des architectures : CNN → Transformers → LLM</li><li>Pourquoi l'attention ? Intuition : “où regarder” dans les données</li><li>Objectifs pédagogiques</li></ul>
Fondements mathématiques des mécanismes d'attention (40 min)	<p>Objectif pédagogique : comprendre l'origine et la rationalité mathématique de l'attention, au-delà de la simple définition calculatoire.</p> <ul style="list-style-type: none"><li>Lecture mathématique : attention comme produit scalaire dynamique</li></ul>

	$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$ <p>mais interprété comme mesure de similarité adaptative entre représentations</p> <ul style="list-style-type: none"> <li>• Produit scalaire dynamique : <math display="block">\langle x, y \rangle_A = x^T A y</math> <p><math>A</math> dépendant du contexte (appris)</p> </li> <li>• Lien avec l'analyse canonique des corrélations (CCA) : <ul style="list-style-type: none"> <li>◦ CCA cherche des projections <math>u, v</math> maximisant <math>\text{corr}(Xu, Yv)</math></li> <li>◦ Attention : projections apprises <math>Q = XW_Q, K = YW_K</math>, similarité calculée localement</li> </ul> </li> </ul> <p><b>Message clé : généralisation de méthodes statistiques classiques pour créer un produit scalaire contextuel et dynamique</b></p>
Attention dans les modèles de langage (LLM) (30 min)	<ul style="list-style-type: none"> <li>• Self-attention et multi-head attention</li> <li>• Architecture Transformer simplifiée</li> <li>• Compréhension de la dépendance globale dans les séquences</li> </ul>
Passage du NLP à la vision (30 min)	<ul style="list-style-type: none"> <li>• Différences texte vs image : structure spatiale vs séquentielle</li> <li>• Tokenisation des images (patches)</li> <li>• Adaptation des mécanismes d'attention à l'espace visuel</li> </ul>
Attention en vision par ordinateur (30 min)	<ul style="list-style-type: none"> <li>• Vision Transformers (ViT)</li> <li>• Attention spatiale et canal (SE, CBAM)</li> <li>• Architectures hybrides CNN + attention</li> </ul>
Session pratique et discussion (60 min)	<ul style="list-style-type: none"> <li>• Implémentation d'un module de self-attention (TensorFlow ou PyTorch)</li> <li>• Application sur des datasets visuels (CIFAR-10, subset ImageNet et des données médicales)</li> <li>• Visualisation et interprétation des cartes d'attention</li> </ul> <p>Discussion :</p> <ul style="list-style-type: none"> <li>• Limites : coût computationnel, robustesse et interprétabilité</li> <li>• Extensions récentes et applications multimodales</li> <li>• Ouverture vers attention contextuelle adaptative</li> </ul>

## Public cible

- Doctorants, chercheurs et ingénieurs en IA
- Enseignants souhaitant introduire l'attention dans leurs cours
- Niveau : intermédiaire

## Références

Dosovitskiy, Alexey, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit and Neil Houlsby. "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale." *ArXiv abs/2010.11929* (2020): n. pag.

Gu, Albert and Tri Dao. "Mamba: Linear-Time Sequence Modeling with Selective State Spaces." *ArXiv abs/2312.00752* (2023): n. pag.

Vaswani Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17). Curran Associates Inc., Red Hook, NY, USA, 6000–6010.

Wang, Y., Wang, W., Li, Y. et al. An attention mechanism module with spatial perception and channel information interaction. *Complex Intell. Syst.* 10, 5427–5444 (2024). <https://doi.org/10.1007/s40747-024-01445-9>