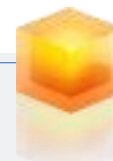


Symposium SFP : Intelligence artificielle et anatomie pathologique

Carrefour Pathologie 2019 – 6 novembre 2019 – Palais des congrès, Paris

Fouille de données, connaissances structurées et apprentissage peu supervisé La clé de l'aide au diagnostic en anatomie pathologie



François-Xavier Frenois, PhD

fx.frenois@canceropole-gso.org



Plateforme Imag'IN de l'IUC : <https://www.imagin.univ-tlse3.fr/ImagIn>



GitHub : <https://github.com/ouatataz>

AUCUN LIEN D'INTÉRÊT



L'intelligence artificielle est née dans les années 1950

Dartmouth Conference, 1956



FAIRE PRODUIRE DES TÂCHES HUMAINES PAR DES MACHINES QUI MIMENT L'ACTIVITÉ DU CERVEAU

Au cours des décennies qui ont suivi, l'IA a été présentée alternativement comme une grande opportunité de développement de nos civilisations, et en même temps très souvent décriée (manque de puissance informatique)

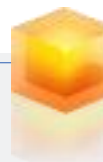


Deux courants de recherche se sont constitués...

- ▶ **IA dite « forte » : General AI**
 - Concevoir des machines possédants tous nos sens et capables de raisonner comme les humains
- ▶ **IA dite « faible » : Narrow AI**
 - Concevoir des machines capables d'aider les humains dans leurs tâches



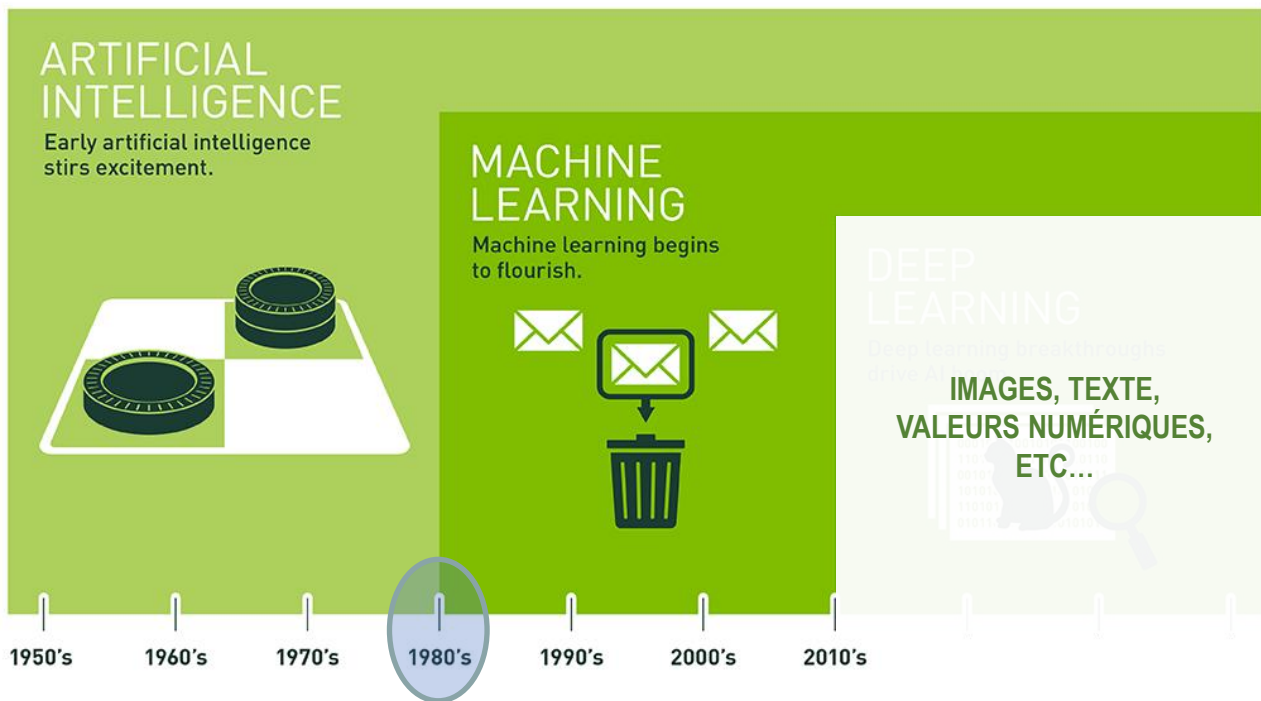
- **Mobilise de nombreuses disciplines** (mathématiques, informatique, sciences cognitives, philosophie) + connaissances spécialisées des domaines auxquels on souhaite l'appliquer (médecine, transports, etc...)
- **Systemes limités dans leurs capacité d'adaptation** : doivent être adaptés pour accomplir des tâches pour lesquelles ils n'ont pas été initialement conçus



L'apprentissage automatique

À la base des technologies mises en œuvre dans l'IA faible

- ▶ **Algorithmes qui raisonnent sur des données : Big Data**
 - *Machine Learning depuis les années 1980 (apprentissage machine)*



ACCUMULATION DE DONNÉES
DIMINUTION DES COÛTS DE STOCKAGE

Ces algorithmes se fondent sur des approches statistiques pour donner aux ordinateurs la capacité d'apprendre à partir de données puis de faire des prédictions sur de nouvelles données

▶ AMÉLIORATION DES PERFORMANCES POUR RÉSOUDRE DES TÂCHES SANS ÊTRE EXPLICITEMENT PROGRAMMÉS POUR CHACUNE



Machine learning : modèle statistique ou probabiliste

Différentes approches algorithmiques souvent utilisées en combinaison pour aboutir à une classification

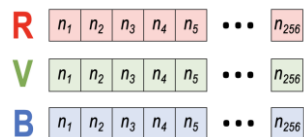
- ▶ **Apprentissage supervisé et non supervisé** : selon les informations disponibles durant la phase d'apprentissage...
 - **Supervisé** : les données d'entrée sont fournies au système avec leur label (i.e. image H&E + Diagnostic associé)
 - **Non supervisé** : pas de label, on cherche à déterminer la structure sous jacentes des données
 - **Autres** : semi-supervisé, par renforcement, transfert learning...

- ▶ **Données possibles** : vecteurs de caractéristiques (n dimensions), graphes, arbres, courbes...

Image RVB 8 bits (256 niveaux de gris / canal)



DOMINANTE VERTE



VECTEURS DE CARACTÉRISTIQUES



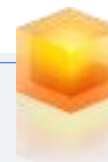
DONNÉES D'ENTRÉ DU CLASSIFIEUR



Image RVB 8 bits (256 niveaux de gris / canal)



DOMINANTE JAUNE



Machine learning : modèle statistique ou probabiliste

Différentes approches algorithmiques souvent utilisées en combinaison pour aboutir à une classification

- ▶ **Apprentissage supervisé et non supervisé** : selon les informations disponibles durant la phase d'apprentissage...
 - **Supervisé** : les données d'entrée sont fournies au système avec leur label (i.e. image H&E + Diagnostic associé)
 - **Non supervisé** : pas de label, on cherche à déterminer la structure sous jacentes des données
 - **Autres** : semi-supervisé, par renforcement, transfert learning...

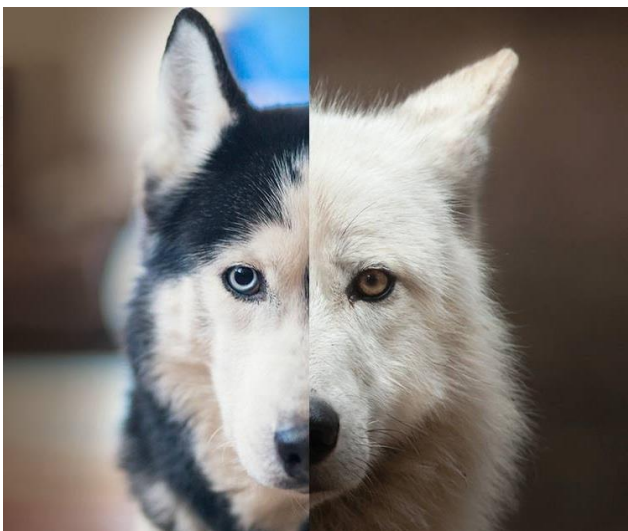
- ▶ **Données possibles** : vecteurs de caractéristiques (n dimensions), graphes, arbres, courbes...

Image RVB 8 bits (256 niveaux de gris / canal)



DOMINANTE VERTE

R	α_1	α_2	α_3	α_4
V	β_1	β_2	β_3	β_4
B	γ_1	γ_2	γ_3	γ_4



α_1	α_2	α_3	...	α_{256}
β_1	β_2	β_3	...	β_{256}
γ_1	γ_2	γ_3	...	γ_{256}

Image RVB 8 bits (256 niveaux de gris / canal)



DOMINANTE JAUNE



Machine learning : modèle statistique ou probabiliste

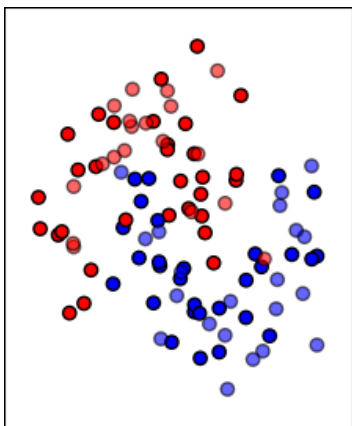
Principaux algorithmes et code de test

▶ Principaux algorithmes :

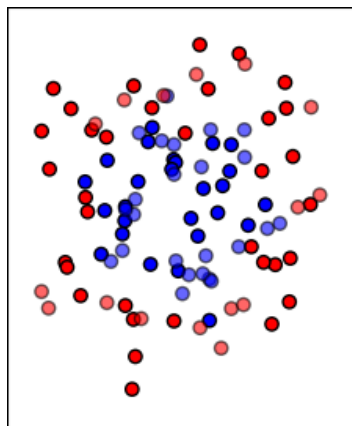
- **Supervisé** : SVM, ARBRES DE DÉCISION (RANDOM FOREST), KNN, PERCEPTRON (RÉSEAU DE NEURONES)...
- **Non supervisé** : K-MEANS...

▶ Python : librairie « scikit-learn »

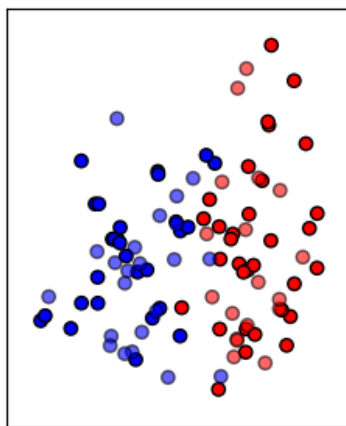
<https://scikit-learn.org/stable/index.html>



DATASET #1
make_moons



DATASET #2
make_circles

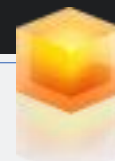


DATASET #3
linearly_separable

```

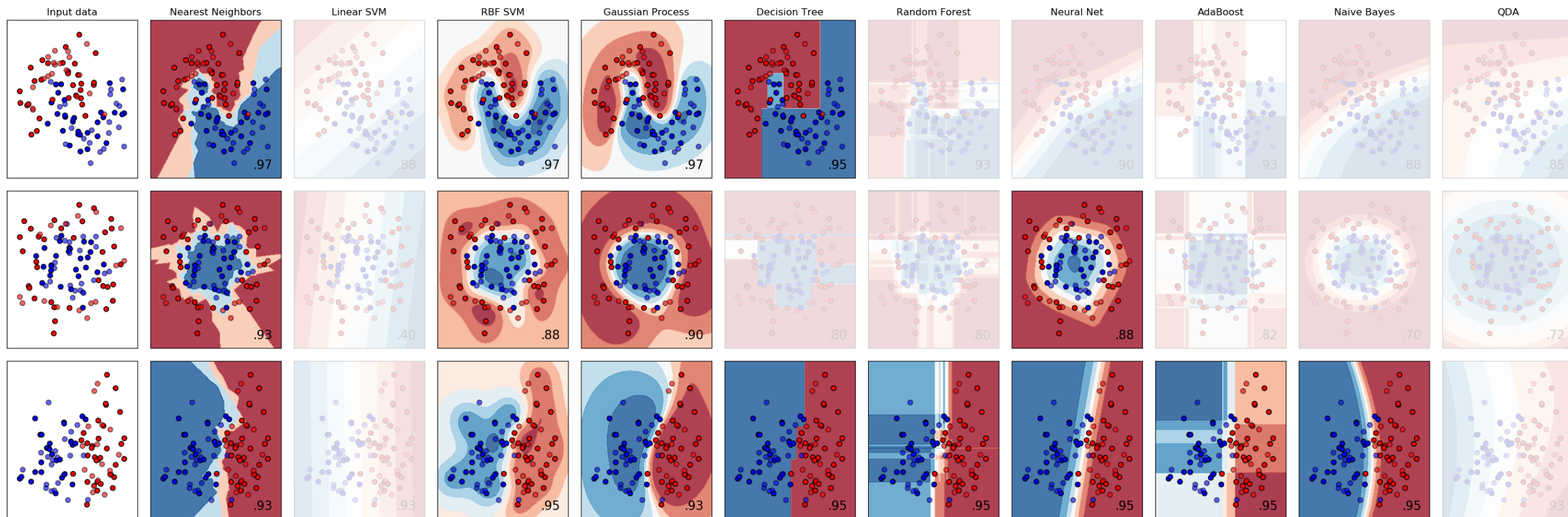
1 # Initialize classifiers
2 names = [
3     "Nearest Neighbors",
4     "Linear SVM",
5     "RBF SVM",
6     "Gaussian Process",
7     "Decision Tree",
8     "Random Forest",
9     "Neural Net",
10    "AdaBoost",
11    "Naive Bayes",
12    "QDA"
13 ]
14 classifiers = [
15     KNeighborsClassifier(3),
16     SVC(kernel="linear", C=0.025),
17     SVC(gamma=2, C=1),
18     GaussianProcessClassifier(1.0 * RBF(1.0)),
19     DecisionTreeClassifier(max_depth=5),
20     RandomForestClassifier(max_depth=5, n_estimators=10, max_features=1),
21     MLPClassifier(alpha=1, max_iter=1000),
22     AdaBoostClassifier(),
23     GaussianNB(),
24     QuadraticDiscriminantAnalysis()
25 ]
26
27 # Generate a random n-class classification problem
28 X, y = make_classification(
29     n_features=2,
30     n_redundant=0,
31     n_informative=2,
32     random_state=1,
33     n_clusters_per_class=1
34 )
35
36 # Initialize datasets
37 datasets = [
38     make_moons(noise=0.3, random_state=0),
39     make_circles(noise=0.2, factor=0.5, random_state=1),
40     linearly_separable
41 ]
42
43 # Iterate over datasets and pre-process
44 # Split datasets into training and test part
45 for ds_cnt, ds in enumerate(datasets):
46     X, y = ds
47     X = StandardScaler().fit_transform(X)
48     X_train, X_test, y_train, y_test = \
49         train_test_split(X, y, test_size=.4, random_state=42)
50
51 # Iterate over classifiers
52 for name, clf in zip(names, classifiers):
53     clf.fit(X_train, y_train)
54     score = clf.score(X_test, y_test)

```



Machine learning : modèle statistique ou probabiliste

Principaux algorithmes et code de test

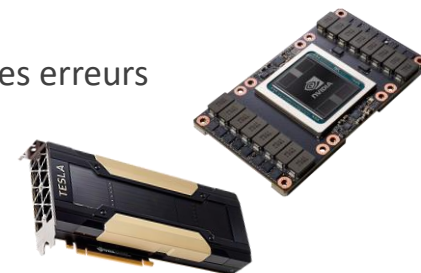


Réseaux de neurones : algorithmes inspirés du fonctionnement du cerveau biologique

Simulent un réseau de neurones organisés en différentes couches, échangeant les uns avec les autres

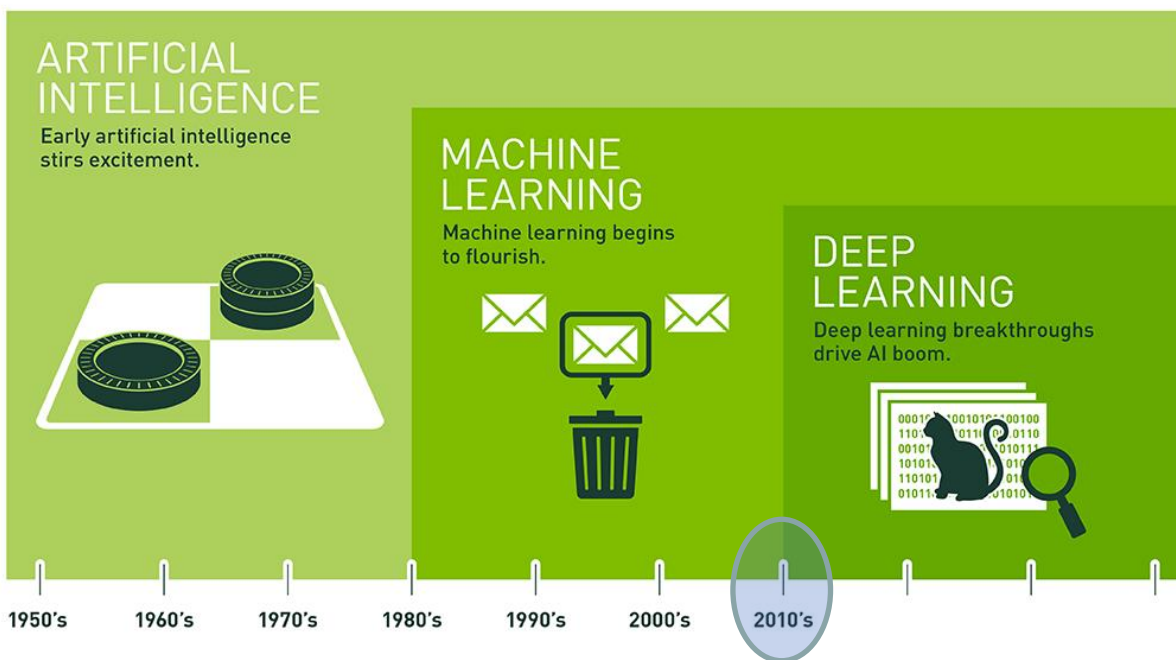
► Réseaux de neurones profonds convolutionnels (CNN) : classification d'images

- **Apprentissage par essais et erreurs** : renforcement des poids synaptiques entre les couches et rétro-propagation des erreurs
- **Très utilisés depuis 2012** : augmentation de la puissance de calcul parallélisable offerte par les cartes GPU...



CARTES GPU
NVIDIA TESLA V100

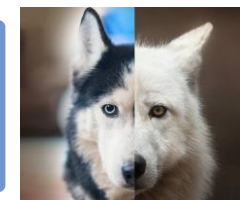
Particularité des réseaux profonds : le nombre de couches est important (niveau d'abstraction très élevé).



CLASSIQUE

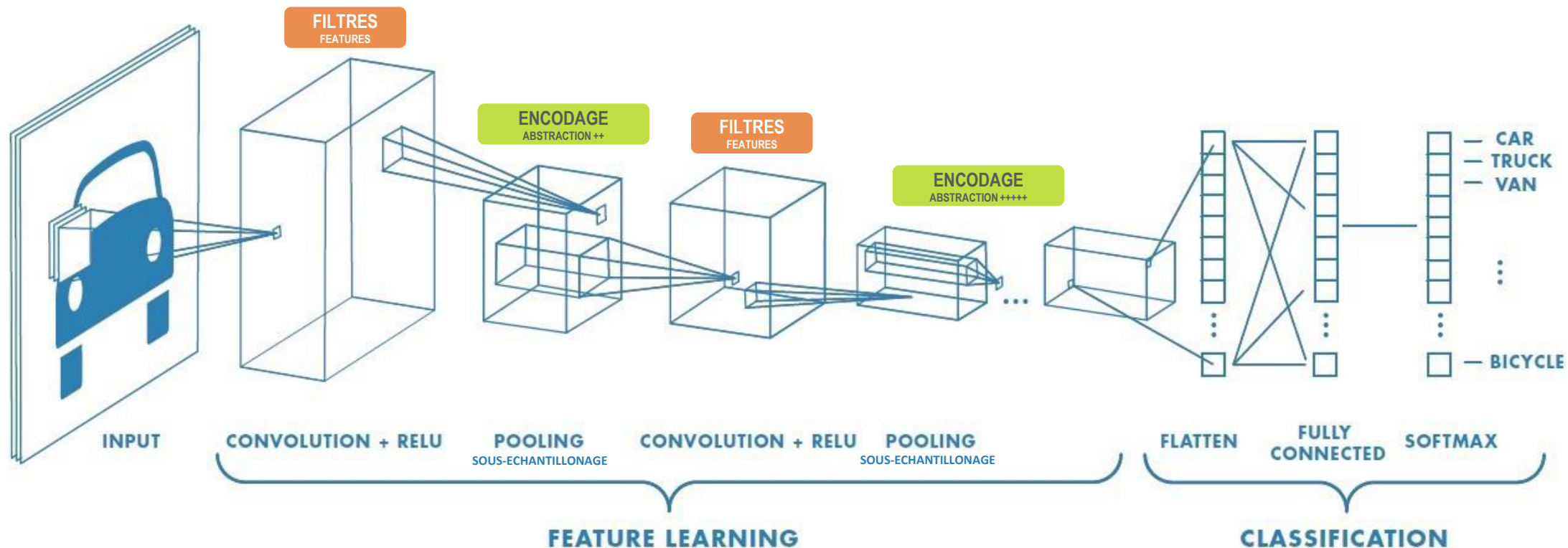


PROFOND



Design d'un CNN profond

Convolution + MaxPooling + Perceptron

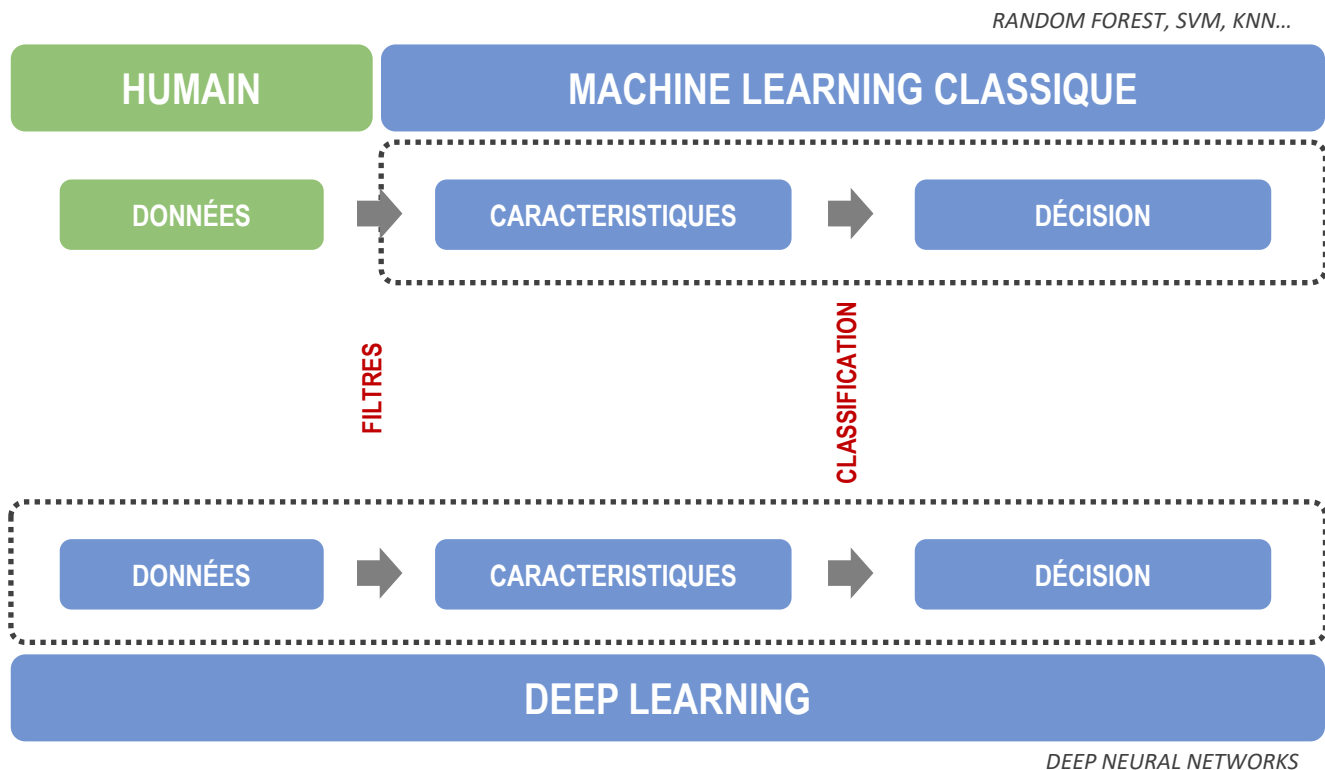


Déterminer automatiquement un espace de caractéristique tout en encodant l'image



Machine learning versus deep learning

L'espace des caractéristiques est connu et défini par l'homme dans le machine learning, alors qu'un réseau de neurones va lui-même déterminer les caractéristiques qui expliquent les données pour permettre la prise de décision



Mise en oeuvre (infrastructures et algorithmes éprouvés)
Caractéristiques connues : hiérarchie, corrélations...



Intervention humaine : objectivité et complexité diminuées
Temps perdu en design de caractéristiques et métriques



Pas d'élaboration de caractéristiques, décisions très complexes
Variété de données : images, langage, etc...



Lourd à mettre en oeuvre
Quantité de données pour l'apprentissage : $10^5 \sim 10^6$
Boîte noire (caractéristiques inconnues), explicabilité



Arnaud Abreu

aabreu@unistra.fr

arnaud.abreu@roche.com

 **GitHub** : <https://github.com/ArnaudAbreu>

Employé de l'Institut Roche



Deep-Learning, principales limitations et mauvaises tendances :

- **Mise en place** : Rassembler, séparer, pré-traiter, labéliser, entraîner et paramétrer, post-traiter (Compétences + Temps)
- **Déploiement** : Méthode statistique...
- **Applicatif** : Souvent atomique, très limité et décevant
- **Explicabilité** : Boîte noire, pixels -> décision sans intermédiaire



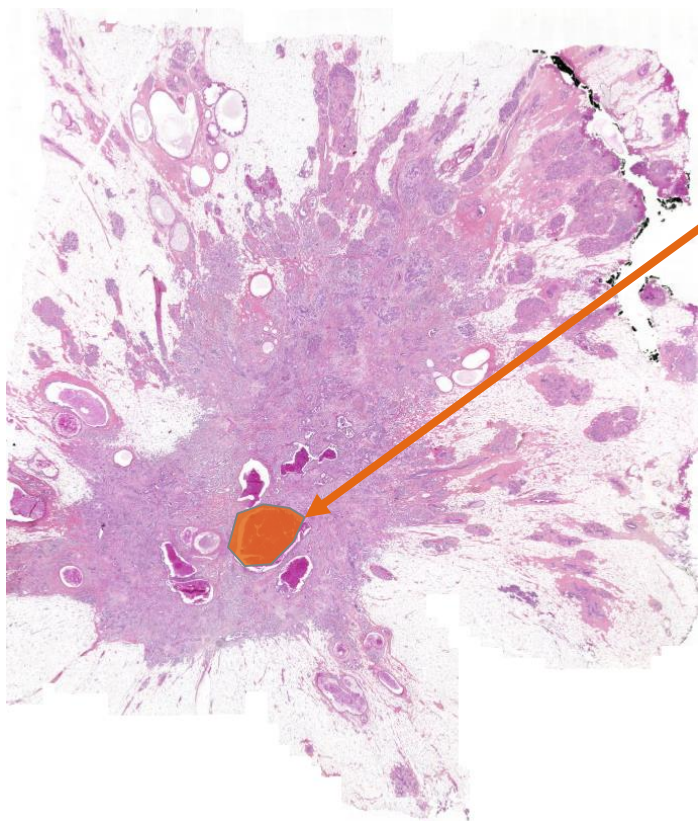
Deep-Learning, principales limitations et mauvaises tendances :

- Mise en place : Rassembler, séparer, pré-traiter **Labéliser**, entraîner et paramétrer, post-traiter (Compétences + Temps)
- Déploiement : Méthode statistique...
- **Applicatif** : Souvent atomique, très limité et décevant → **Exhaustivité**
- **Explicabilité** : Boîte noire, pixels -> décision sans intermédiaire



Exhaustivité et explicabilité de l'analyse : le langage du pathologiste

Formalisme rigoureux de la structuration de la connaissance

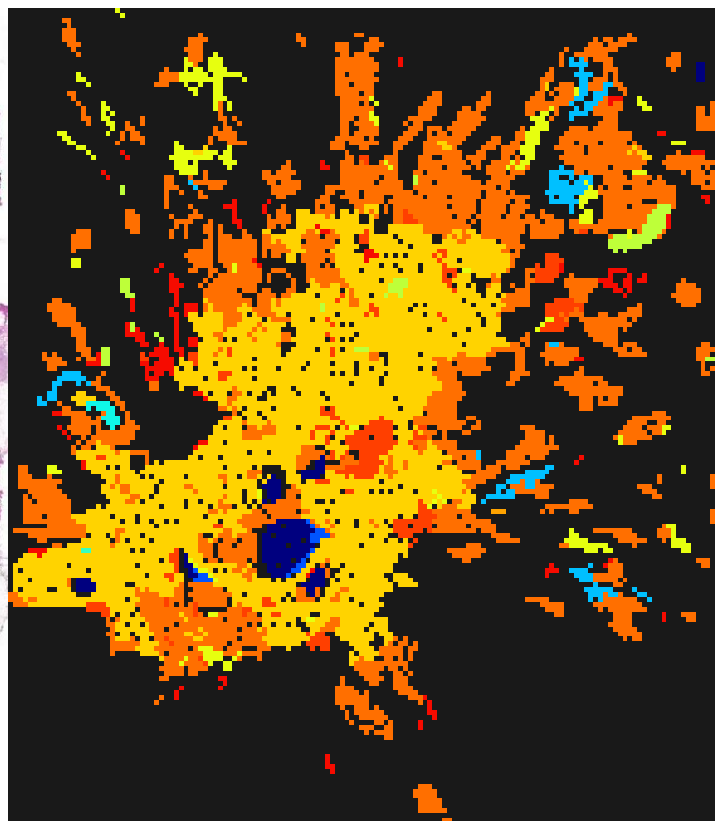
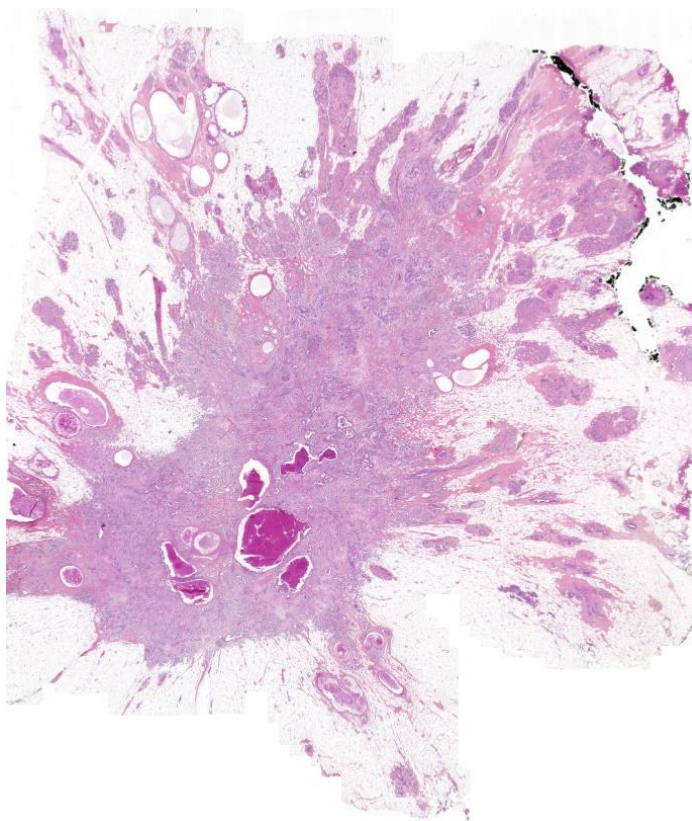


Structure d'intérêt



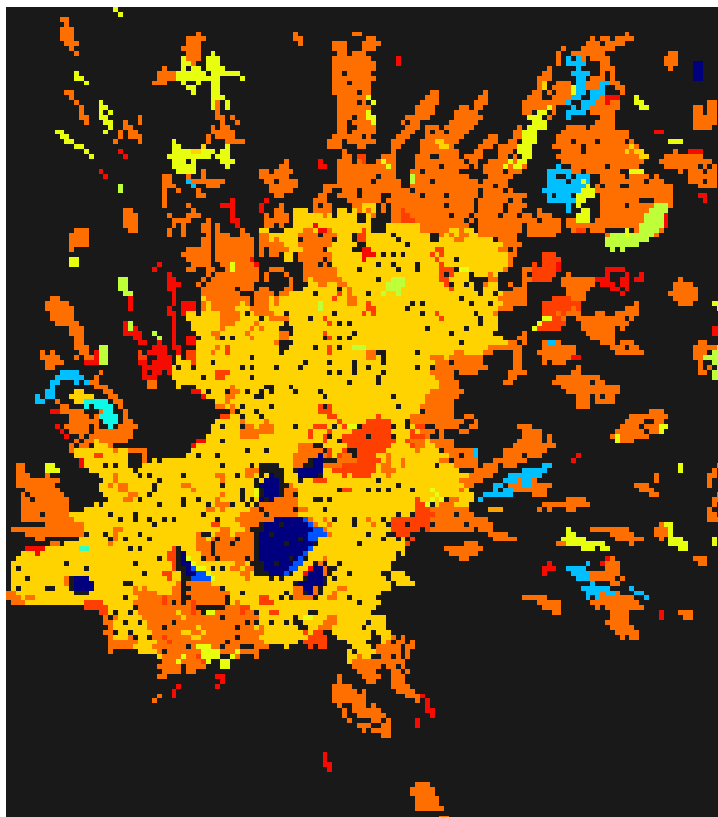
Exhaustivité et explicabilité de l'analyse : le langage du pathologiste

Formalisme rigoureux de la structuration de la connaissance



Exhaustivité et explicabilité de l'analyse : le langage du pathologiste

Formalisme rigoureux de la structuration de la connaissance



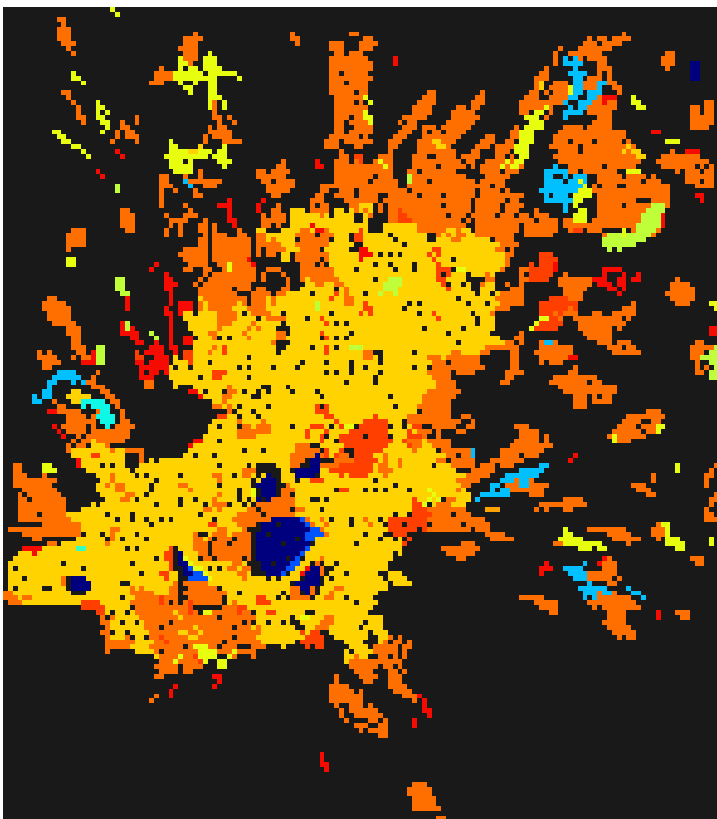
Identification

Segmentation,
classification



Exhaustivité et explicabilité de l'analyse : le langage du pathologiste

Formalisme rigoureux de la structuration de la connaissance

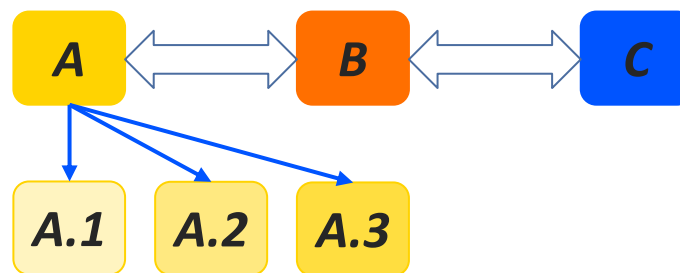


Identification

Segmentation,
classification

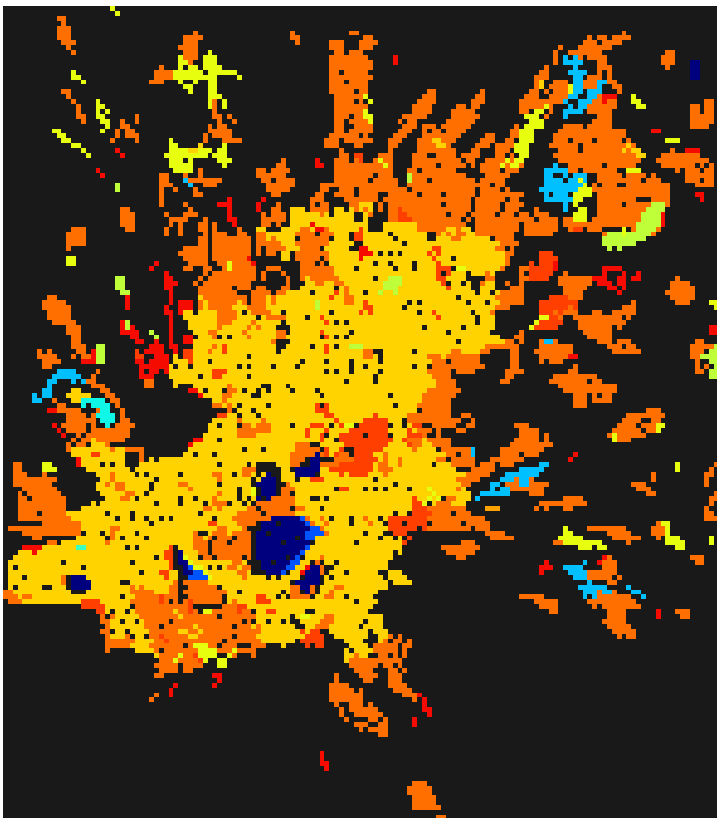
Hiérarchie conceptuelle

Sous-type des structures



Exhaustivité et explicabilité de l'analyse : le langage du pathologiste

Formalisme rigoureux de la structuration de la connaissance



Identification

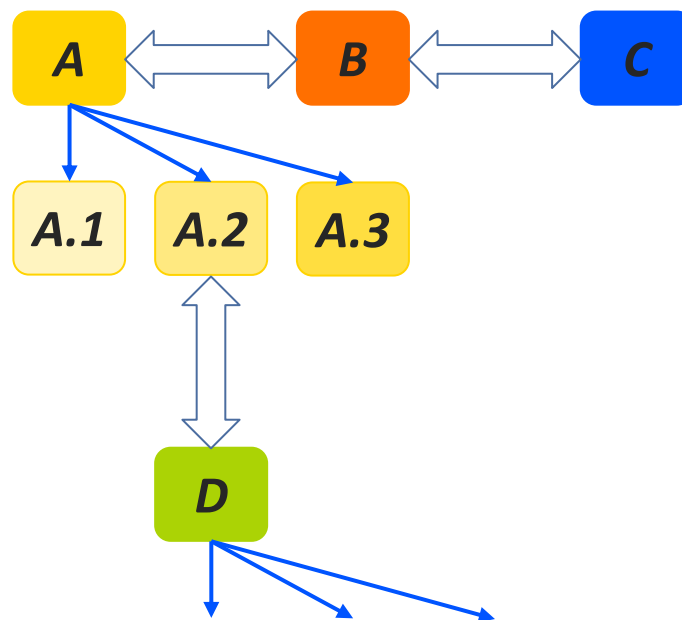
Segmentation,
classification

Hiérarchie conceptuelle

Sous-type des structures

Hiérarchie spatiale

Concepts visibles à plus
forte résolution

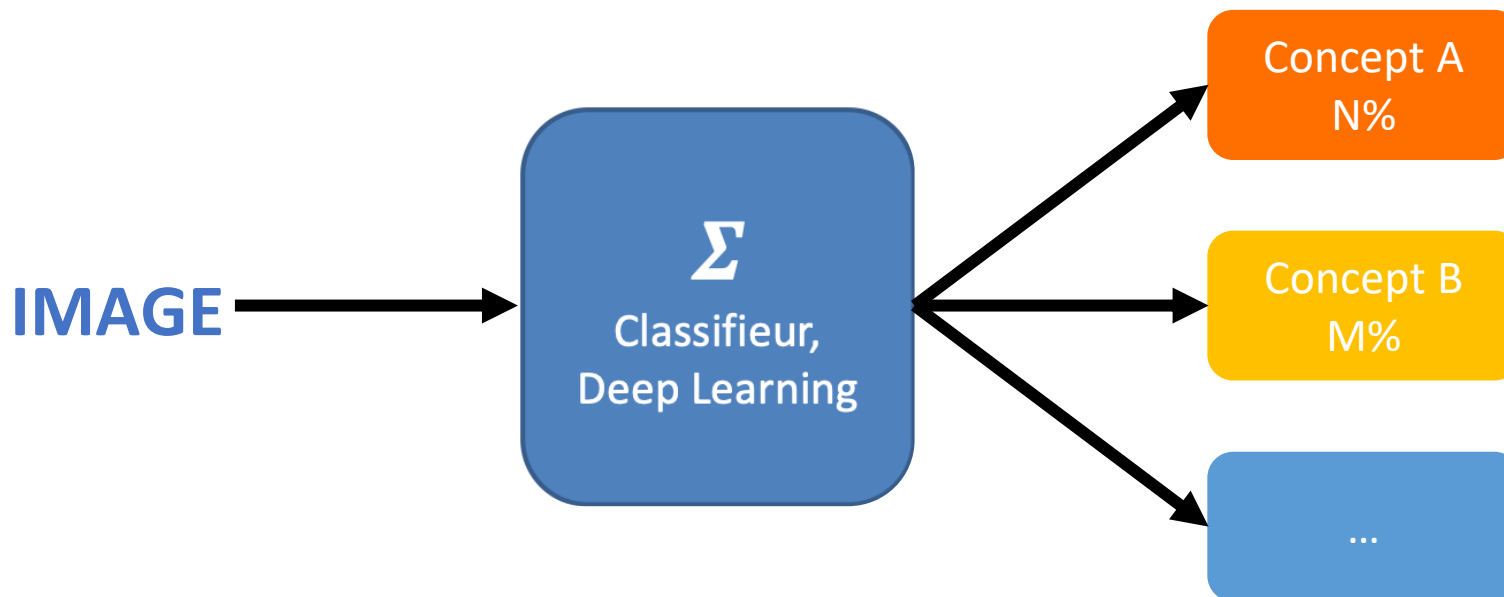


Décision obtenue par déduction



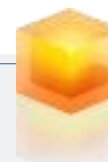
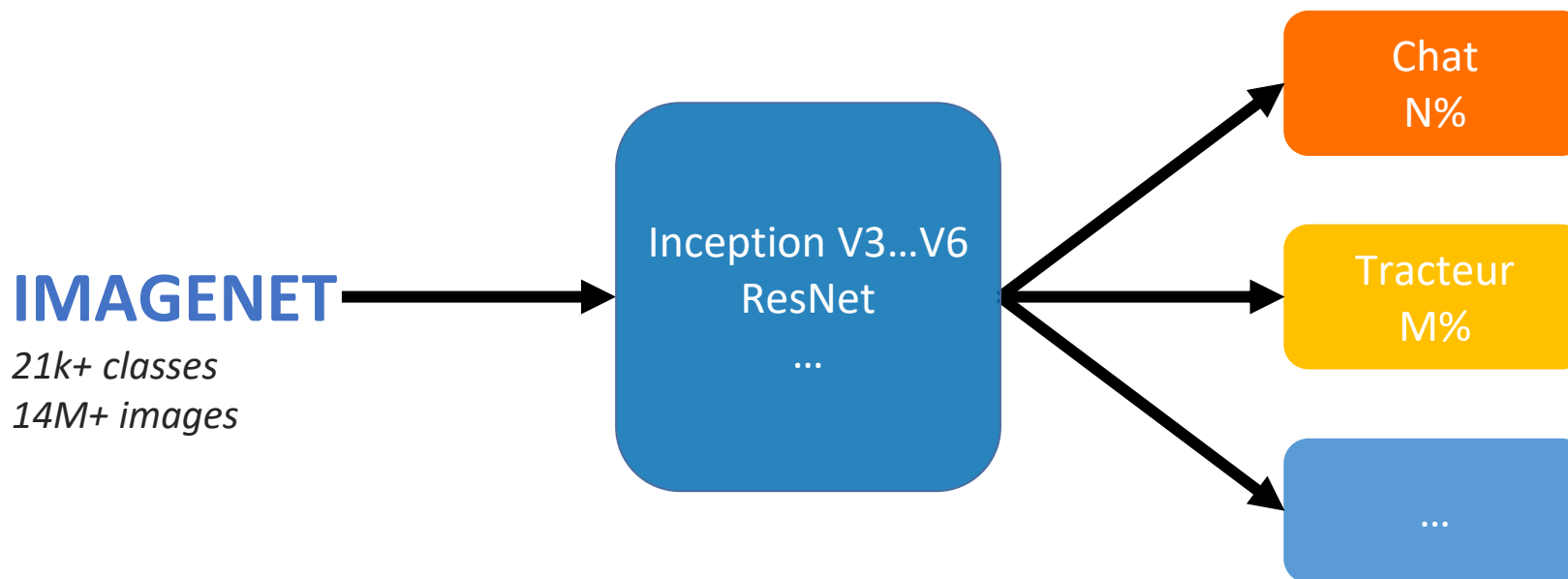
Approche machine-learning, le problème de l'annotation

Un grand nombre de classes, Une grande variabilité, des experts qui ont autre chose à faire...



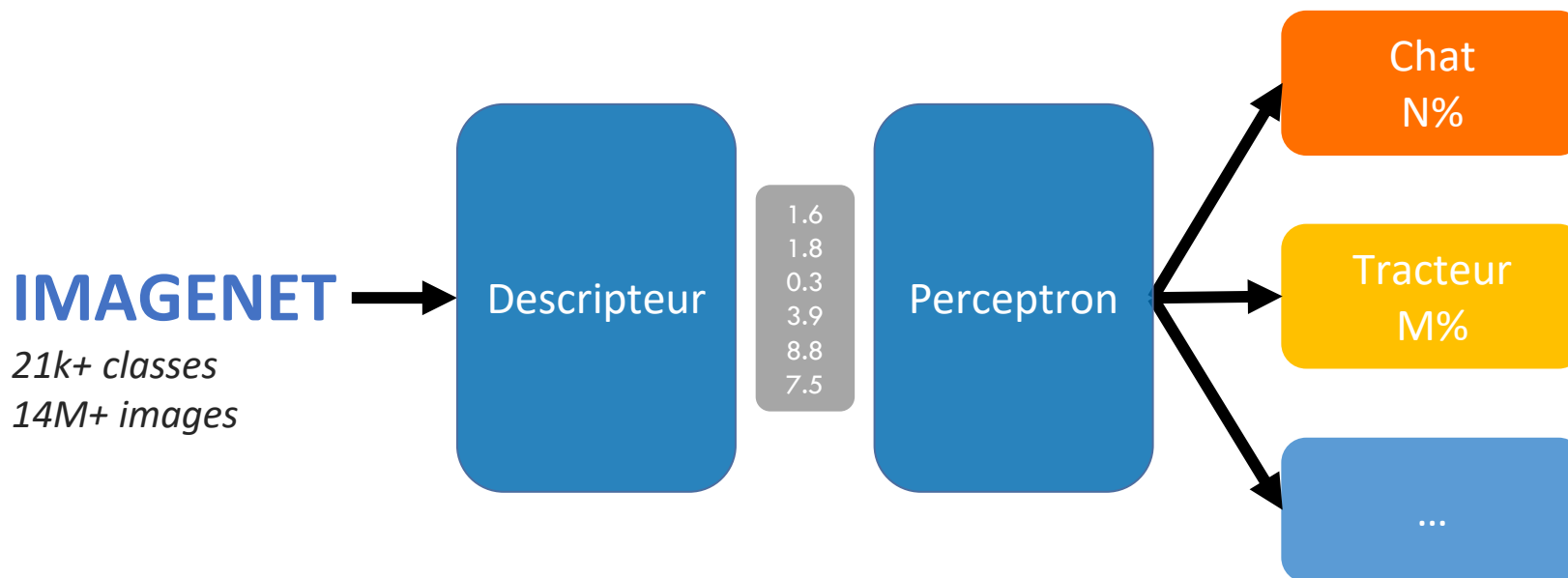
Approche machine-learning, le problème de l'annotation

Transfert de connaissances, utilisation de réseaux entraînés sur des tâches très générales



Approche machine-learning, le problème de l'annotation

Transfert de connaissances, utilisation de réseaux entraînés sur des tâches très générales

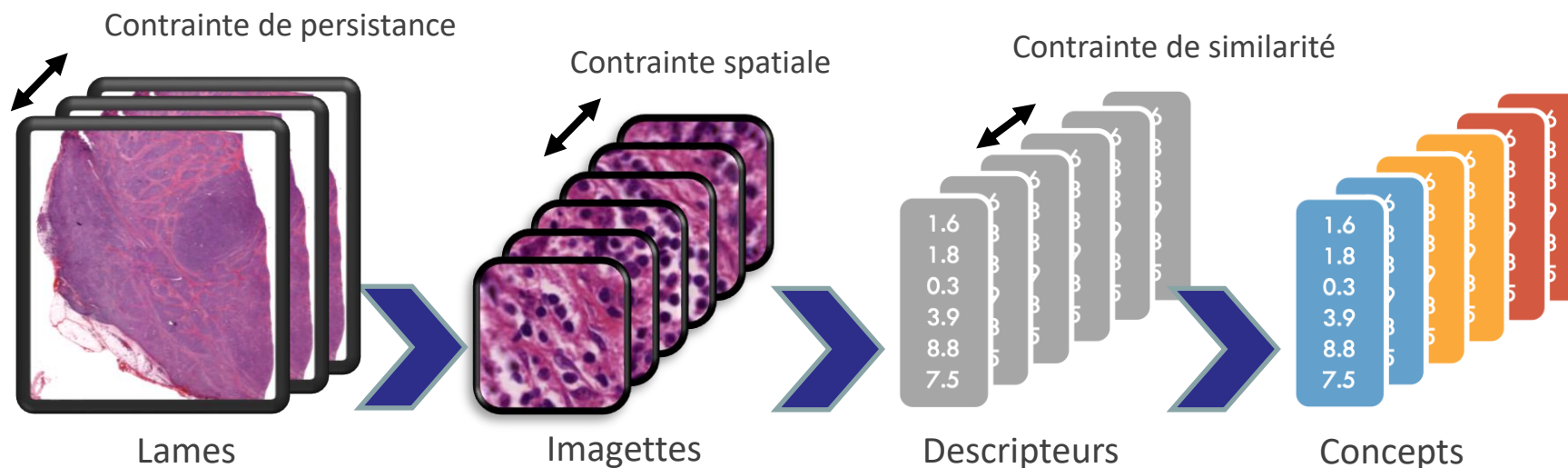


Approche machine-learning, le problème de l'annotation

Transfert de connaissances, clustering dans l'espace des caractéristiques

► Clustering sous certaines contraintes

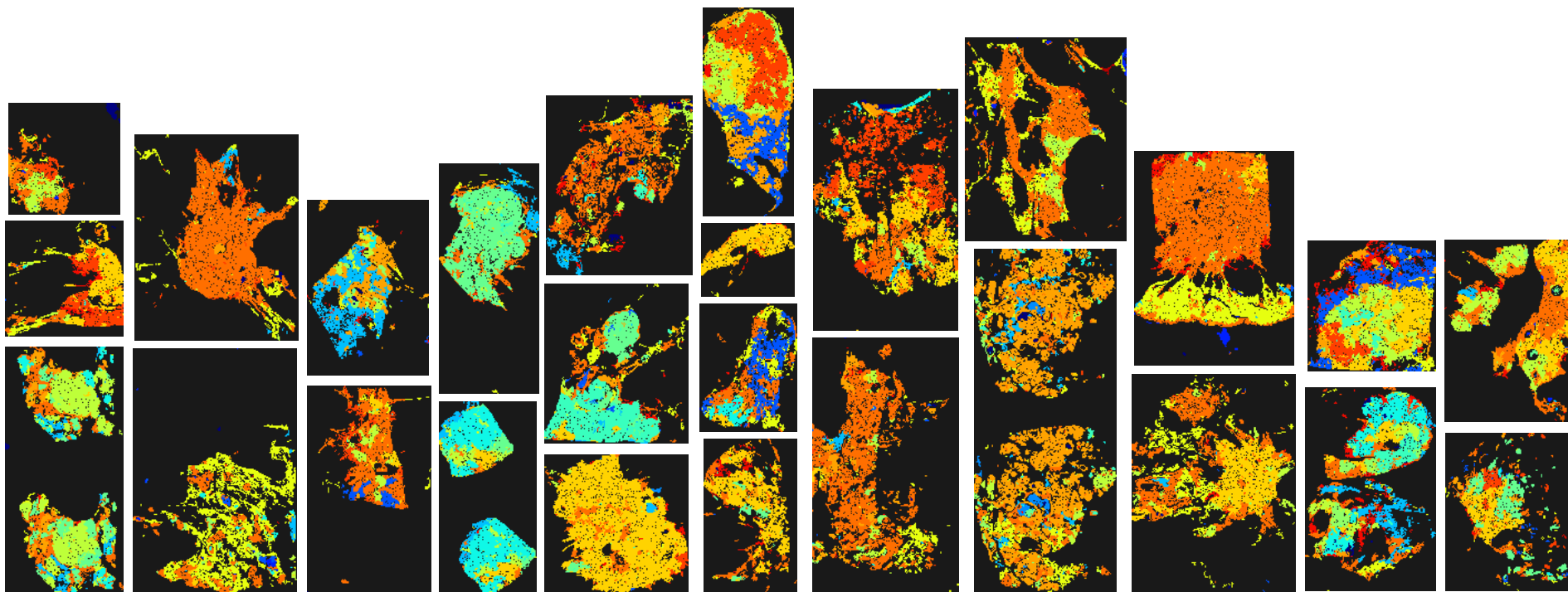
- **Contrainte spatiale de transition lente** : Le concept d'un pixel est sûrement le même que celui de son voisin
- **Contrainte de persistance du motif** : Un concept est observable chez plusieurs patients
- **Contrainte de similarité** : Deux segments d'un même concept doivent avoir une représentation similaire



Applications, perspectives

Usage sur de très grands datasets

- **Sélection de zones d'intérêt** : exploration de lame optimisée
- **Établissement de sous-concepts pertinents** : Émergence de nouveaux biomarqueurs
- **Requête de bases d'images par mots-clefs** : Études cliniques et regroupement de patients facilités



MERCI POUR VOTRE ATTENTION



INSTITUT UNIVERSITAIRE
DU CANCER DE TOULOUSE

Département de Pathologie (CHU Toulouse)

- Pr Pierre Brousset
- Dr François-Xavier Frenois
- Mr Arnaud Abreu
- Dr Camille Franchet
- Dr Camille Laurent
- Me Myriam Marty
- Mr Gaël Gascoin
- ... et l'ensemble du personnel



IMAG'IN

Plateforme Imag'IN de l'IUC

- Dr François-Xavier Frenois
- Dr Nathalie Van Acker
- Me Eveline Caussat
- Me Stéphanie Grenard
- Me Gabrielle Perez



iCUBE

Laboratoire iCube (Strasbourg)

- Dr Cédric Wemmert
- Dr Benoit Naegel



