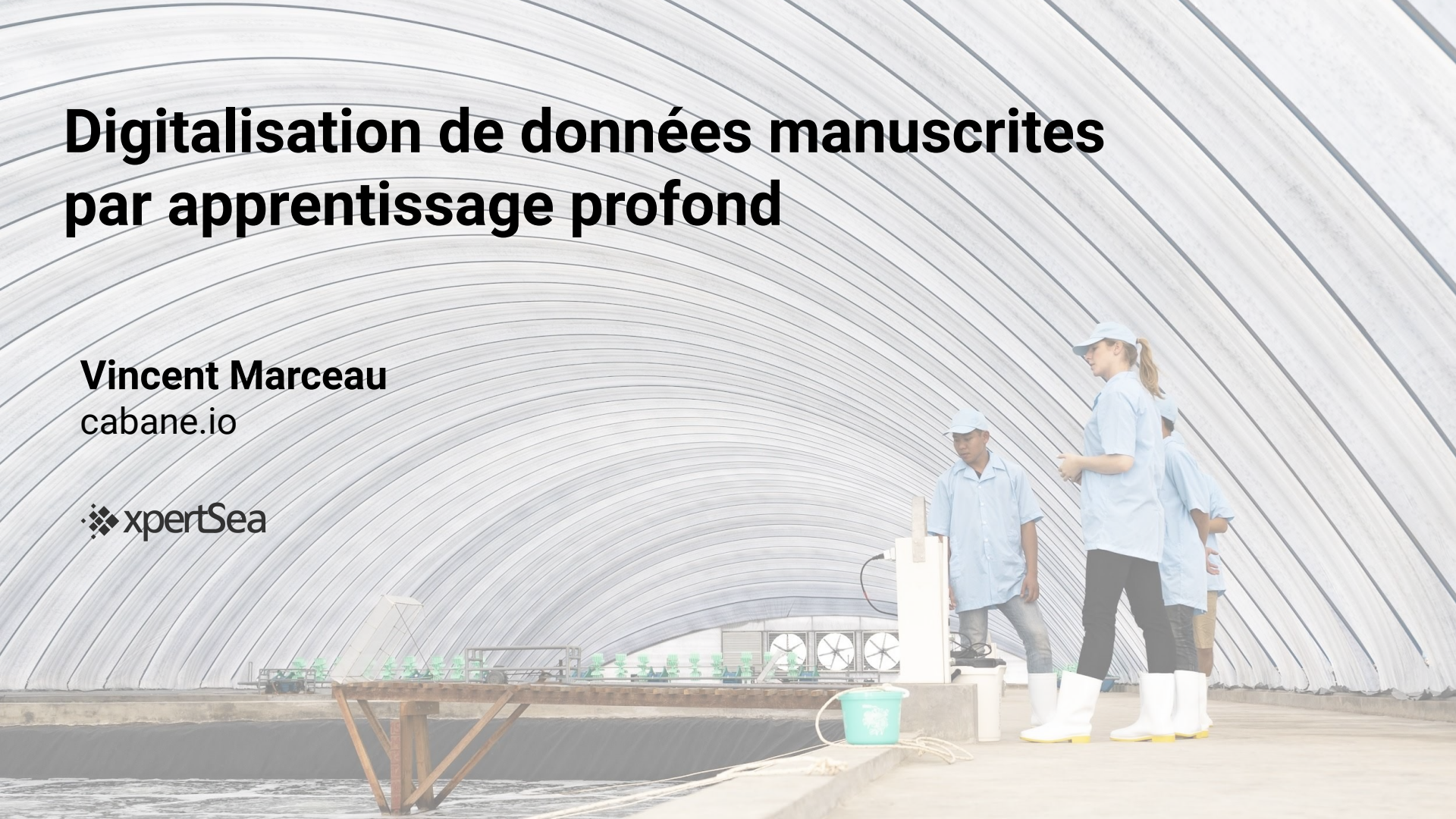


Digitalisation de données manuscrites par apprentissage profond

Vincent Marceau
cabane.io





Digitalisation de données

Pourquoi?

Faciliter l'entrée de données sur notre plateforme.

Comment?

Digitaliser des logs manuscrits des fermiers.

Quoi?

Développer un modèle de reconnaissance de texte manuscrit.

Sample Details			Water Analysis										Ammonia		Microbiology	
Site	pH	Salinity ppt	Ca ppm	Mg ppm	Hardness ppm	CO ₂	HCO ₃	Total	NH ₄	NO ₂	NO ₃	Micro	Yeast	NO ₂	NO ₃	
11	7.6	25	476	1045	4760	NH	150	150	NH	0.91		30	100			
12	7.5	25	312	986	4370	NH	208	208	0.001	1.09		NH	80			
13	7.7	25	272	959	4380	NH	182	182	0.001	1.16		NH	100			
14	7.9	26	294	1016	4785	NH	210	210	NH	1.16		NH	90			
15	7.8	27	314	915	4080	NH	196	196	NH	1.11		NH	100			
16	7.7	27	296	1055	4660	NH	174	174	NH	1.09		NH	100			
17	8.1	25	284	1119	4890	16	176	192	NH	1.18		NH	80			
18	8.2	28	272	942	4150	20	154	174	NH	1.20		NH	100			
19	7.9	28	310	986	4380	NH	208	208	NH	1.20		NH	100			

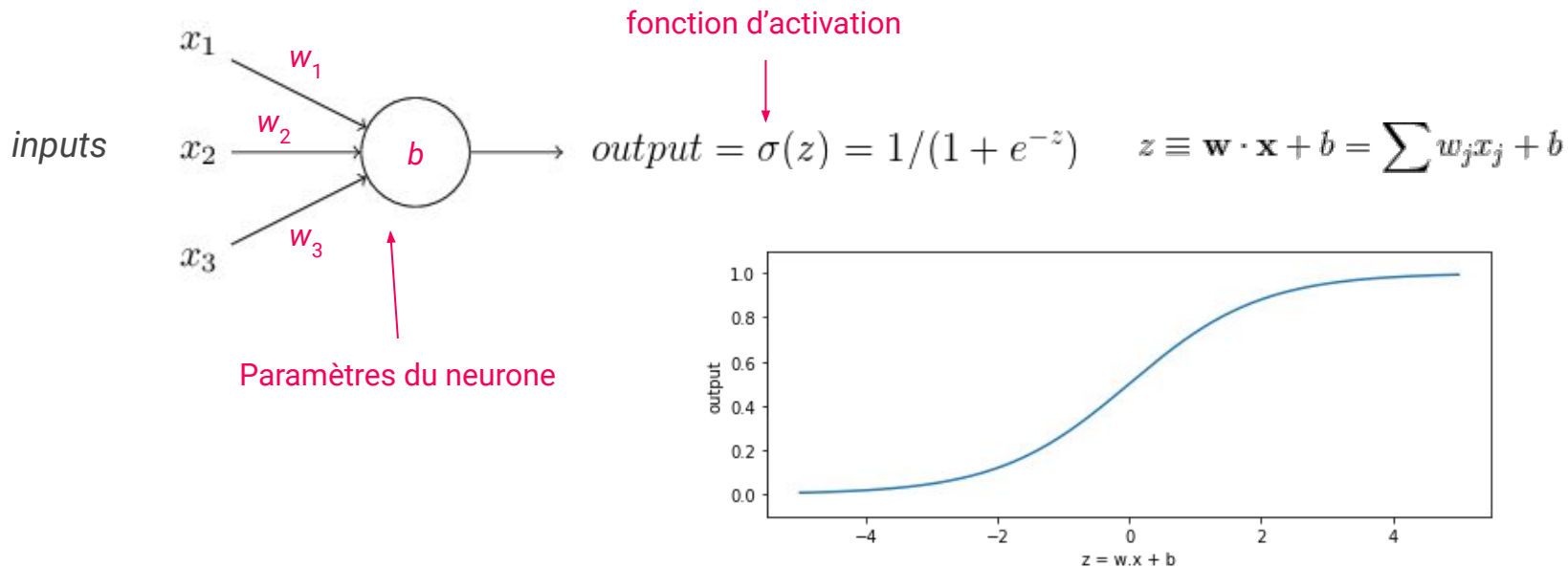
Flial 1 : Test count nauplii								
Date	Master-Hand	Actual	Machine-QB			Machine-BD		
			Count-1	Count-2	Count-3	Count-1	Count-2	Count-3
13/06/18	20000	20000						
13/06/18	30000	30000						
13/06/18	40000	40000				18.616	18.816	19.068
13/06/18	50000	50000				20.227	20.837	21.044
13/06/18	60000	60000				25.950	27.311	29.839
13/06/18	70000	70000				31.933	34.861	37.206
13/06/18	80000	80000				40.830	43.611	49.644
13/06/18	90000	90000				53.091	58.237	60.443
13/06/18	100000	100000				68.322	69.590	65.227
						71.000	69.590	70.972
						80.67	77.6	76.936





Réseaux de neurones: récapitulatif

Le neurone sigmoïde

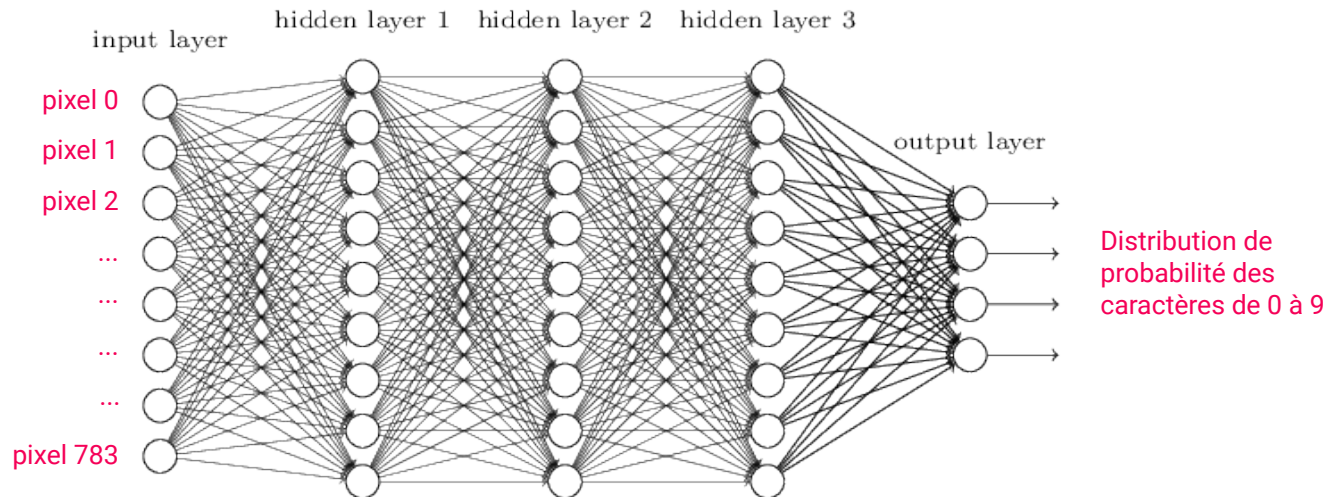
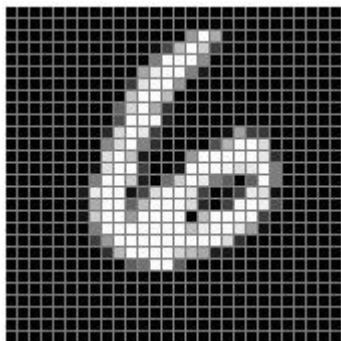




Réseaux de neurones: récapitulatif

Réseau de neurones “fully connected”

Image d'entrée (MNIST)



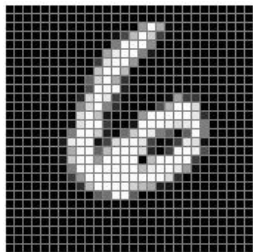


Réseaux de neurones: récapitulatif

Entraînement par descente du gradient

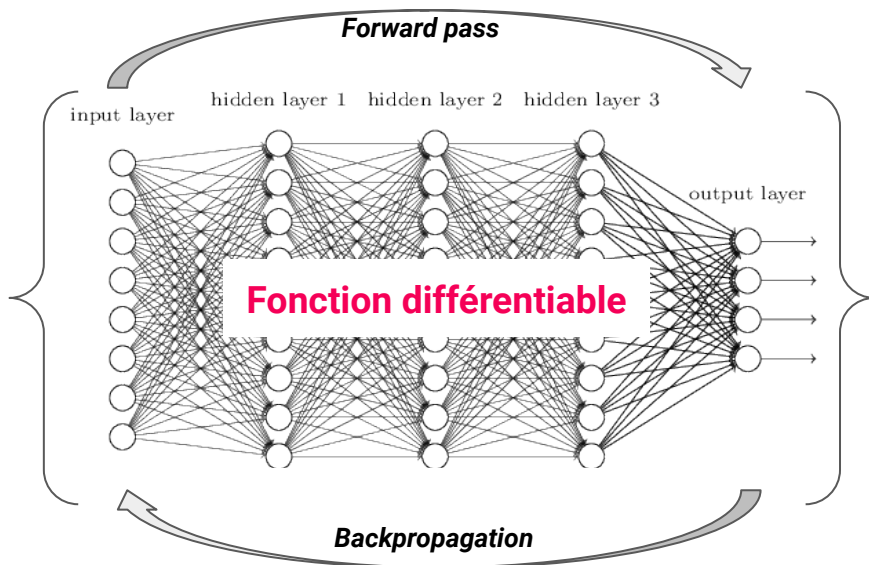
Données d'entraînement

Image annotée:



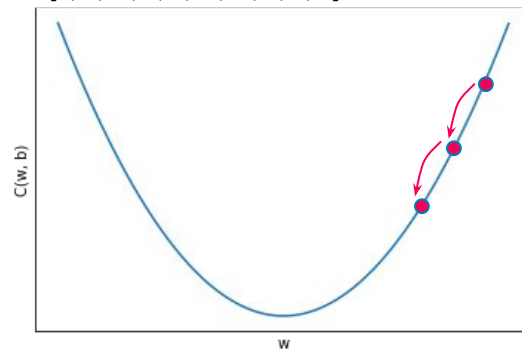
Vérité:

[0, 0, 0, 0, 0, 0, 1, 0, 0, 0]

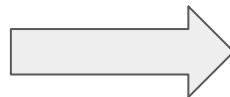
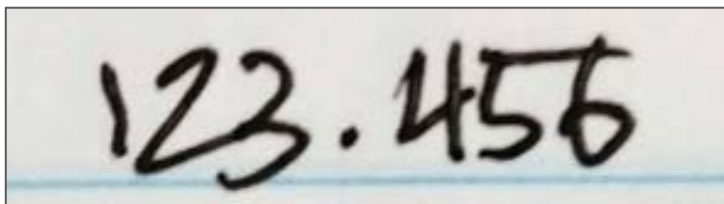


Prédictions et fonction de coût:

[0, 0, .2, 0, 0, .3, .5, 0, 0, 0]



Séquences manuscrites: formulation du problème



“123.456”

Chiffres 0-9 + séparateur décimal

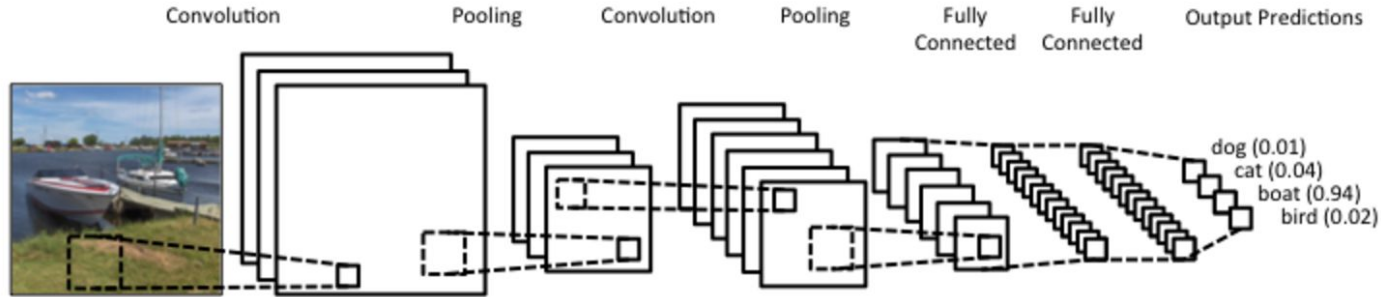
Partie entière et décimale de longueur inconnue

Largeur et position des caractères inconnues

Superposition possible entre les caractères

Approche sélectionnée: Modèle de type **CRNN** (arXiv:1507.05717)

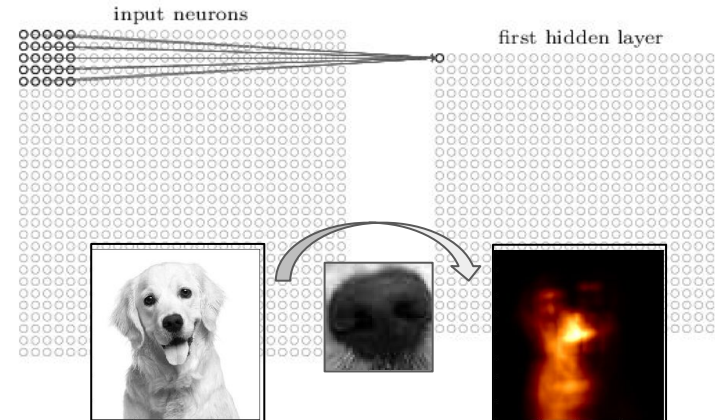
Ingrédient 1: Réseaux de neurones convolutionnels (CNNs)



Pour les problèmes de vision numérique.

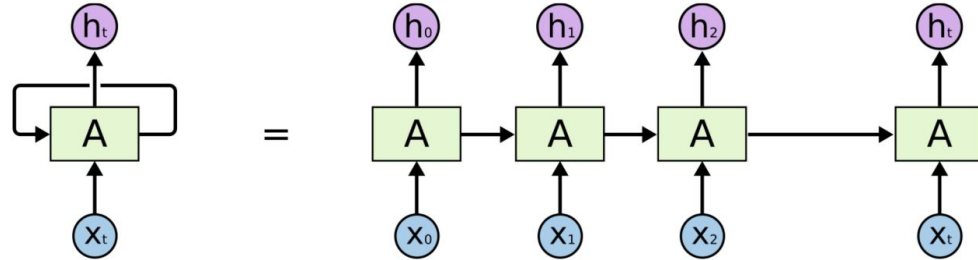
Utilise la cohérence spatiale pour réduire les connexions.

Apprentissage de *kernels* de convolution.





Ingrédient 2: Réseaux de neurones récurrents (RNNs)



Pour les problèmes de prédiction de séquences.

Peut être vu comme une succession de copies du même réseau, chacun passant un “message” à son successeur.

État interne pour l'apprentissage de corrélations temporelles



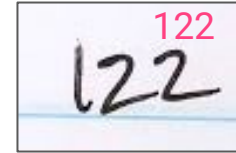
Ingrédient 3: Transcription CTC

Méthode d'encodage / décodage de séquences

Caractère immatériel de séparation “-” pour distinguer les caractères identiques adjacents.

Décodage en 2 étapes:

1. Combinaison des caractères identiques adjacents.
2. Retrait des “-” pour obtenir le mot final.



Exemple de décodage:

“---1111--222-22222”

↓
“-1-2-2-”

↓
“122”



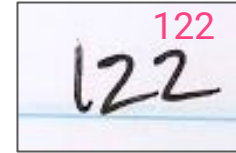
Ingrédient 3: Transcription CTC

Méthode d'encodage / décodage de séquences

Caractère immatériel de séparation “-” pour distinguer les caractères identiques adjacents.

Décodage en 2 étapes:

1. Combinaison des caractères identiques adjacents.
2. Retrait des “-” pour obtenir le mot final.



Encodages valides:

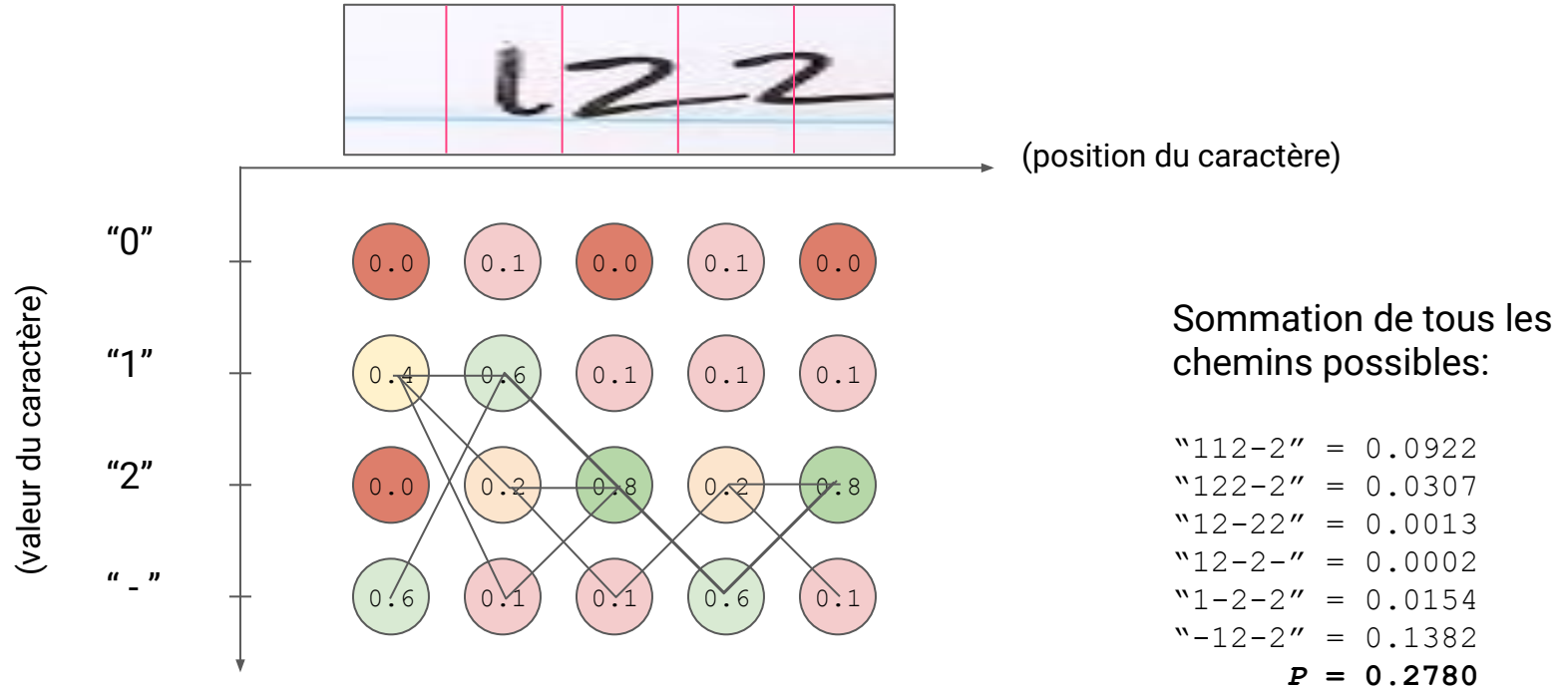
“12-2”, “1122-22”, “--1--2---2”

Encodages invalides:

“122”, “1-22”, “112222”



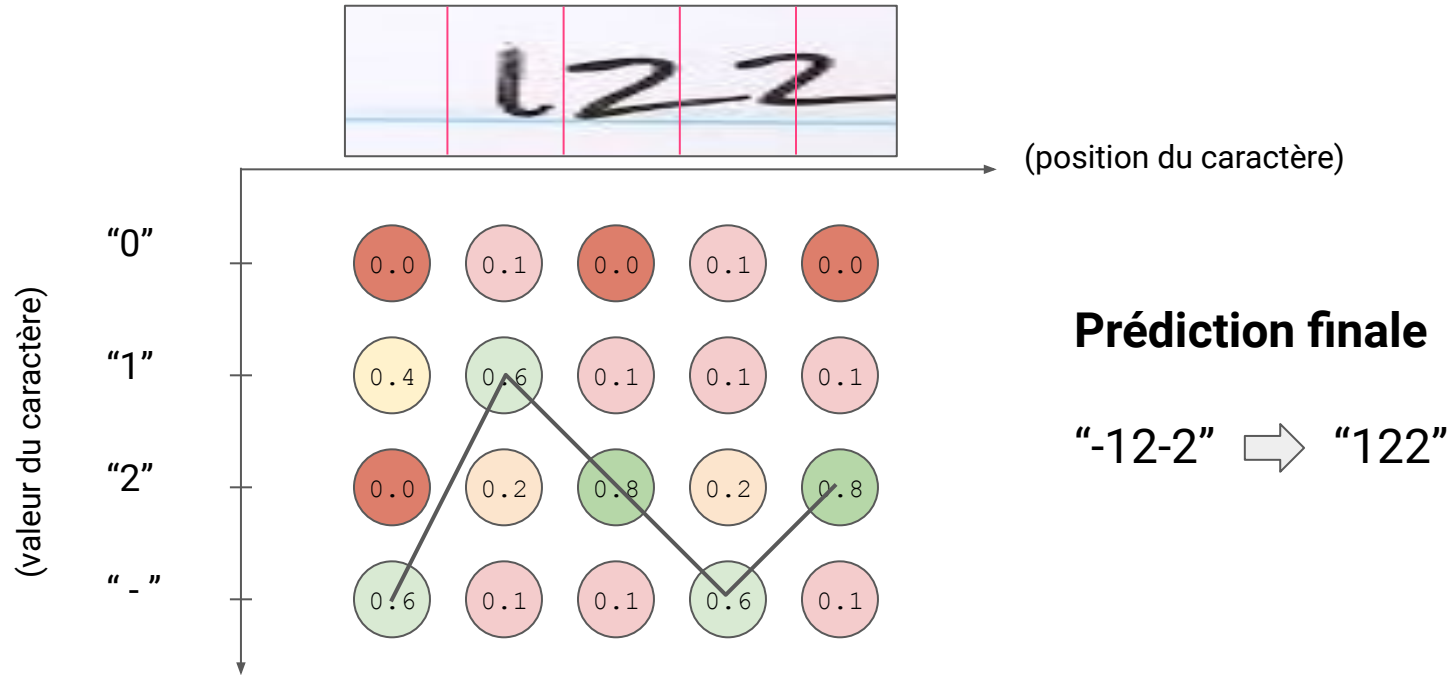
CTC: Calcul des probabilités et fonction de coût



Prédiction = distribution de probabilités



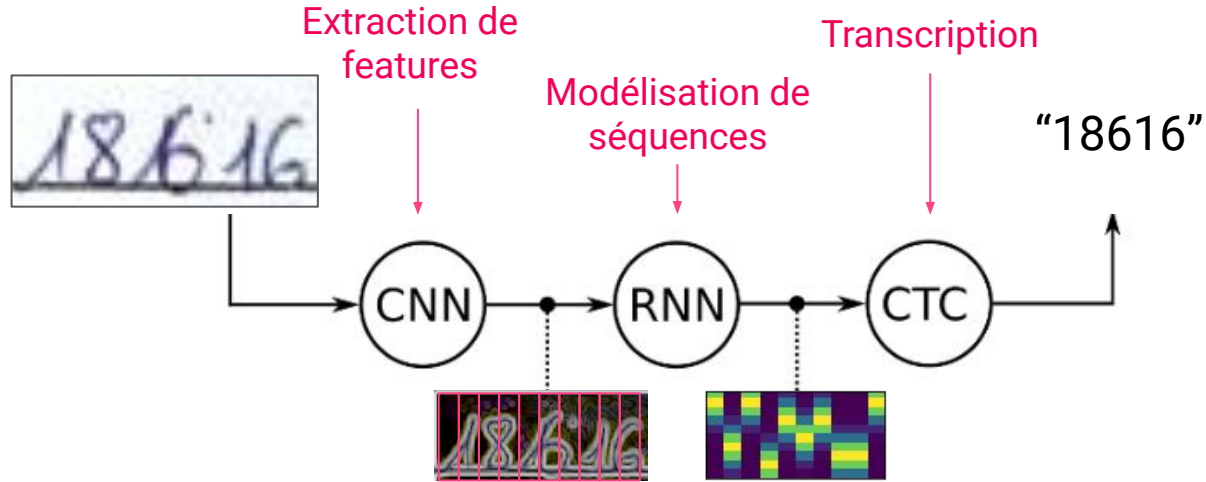
CTC: Calcul des probabilités et fonction de coût



Prédiction = distribution de probabilités



1 + 2 + 3: Architecture CRNN



Modèle entraînable *end-to-end*

Pas nécessaire de segmenter les caractères

Facile de quantifier l'incertitude de la prédiction



Entraînement et données synthétiques

19-8-2049	19	8	2049
ID	3566c3b9-f1ac-4637-8df7-a9e3839a0d61		
Votre nom	Vincent Séguin		

Morning (AM)						
Pond name	Daily feed (kg)	DO (mg / L)	Temperature (°C)	Salinity (g / L)	pH	(m)
1 D27TJ6	1310	5.5	37.4	41	6.60	42
2 D27TJ6	1310	5.5	37.4	41	6.60	42
3 D5TJ3	733	12.69	35.0	90	6.00	19
4 D5TJ3	733	12.69	35.0	90	6.00	19
5 D13TJ9	817	6.71	33.01	52	9.5	27

1. Formulaires "pré-annotés"
2. Génération d'images synthétiques (avec MNIST)

Images réelles



Images synthétiques



Merci



 xpertsea.com

 github.com/vmarceau

 vincent@xpertsea.com

