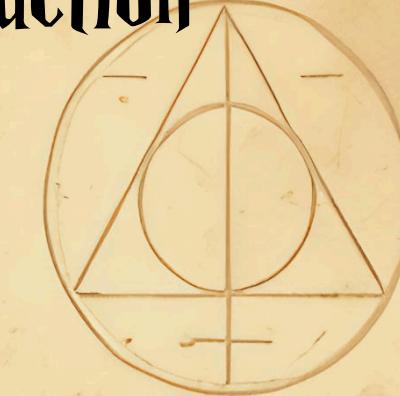


Attaques d'IA - Livre d'introduction



De Sckathapschal Gorpheus Quantifilius Artificewick des Vents

Property of
le magicien quantique

Formule pour cook un challenge du 404

1. Ajouter un chat

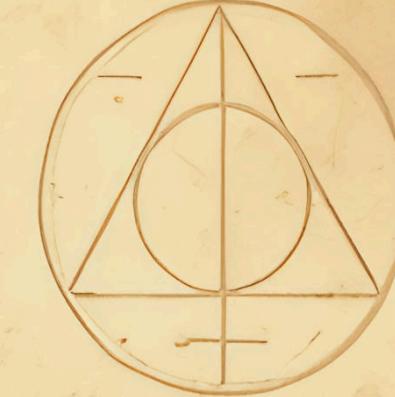


2. Remuer 3 fois dans le sens horaire



3. Ajouter 16 gouttes de potion de Babillage

4. Couper 2 têtes d'hydre, et mélanger le tout



Property of
le magicien quantique

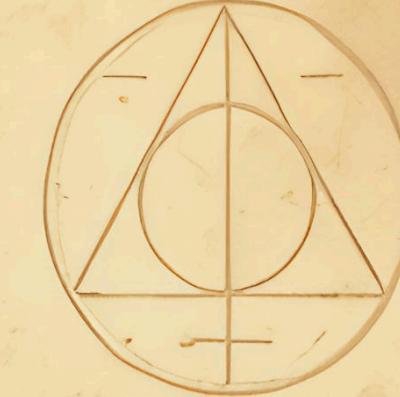
Formule pour cook un challenge du 404

1. Ajouter un chat → une théière

2. Remuer 3 fois dans le sens horaire

3. Ajouter 16 gouttes de potion de Babillage

4. Couper 2 têtes d'hydre, et mélanger le tout



PAGE DE PUB: MAIS QU'EST CE QUE L'IA????

Les débuts (années 1950-1960)

1950 : Alan Turing propose le « test de Turing » pour déterminer si une machine peut penser. 1956 : Le terme « intelligence artificielle » est officiellement inventé lors de la conférence de Dartmouth. Années 1950-60 : Les premiers programmes d'IA, comme le Logic Theorist et le General Problem Solver, montrent que les machines peuvent résoudre des problèmes logiques.

Le premier hiver de l'IA (1970-début des années 1980)

Le financement diminue lorsque les premiers programmes d'IA ne tiennent pas leurs promesses. Les limites de la puissance de calcul deviennent évidentes.

Systèmes experts (années 1980)

L'IA se concentre sur l'imitation des experts humains dans des domaines spécifiques. Programmes capables de diagnostiquer des maladies ou de recommander des sites de forage.

L'essor de l'apprentissage automatique (années 1990-2000)

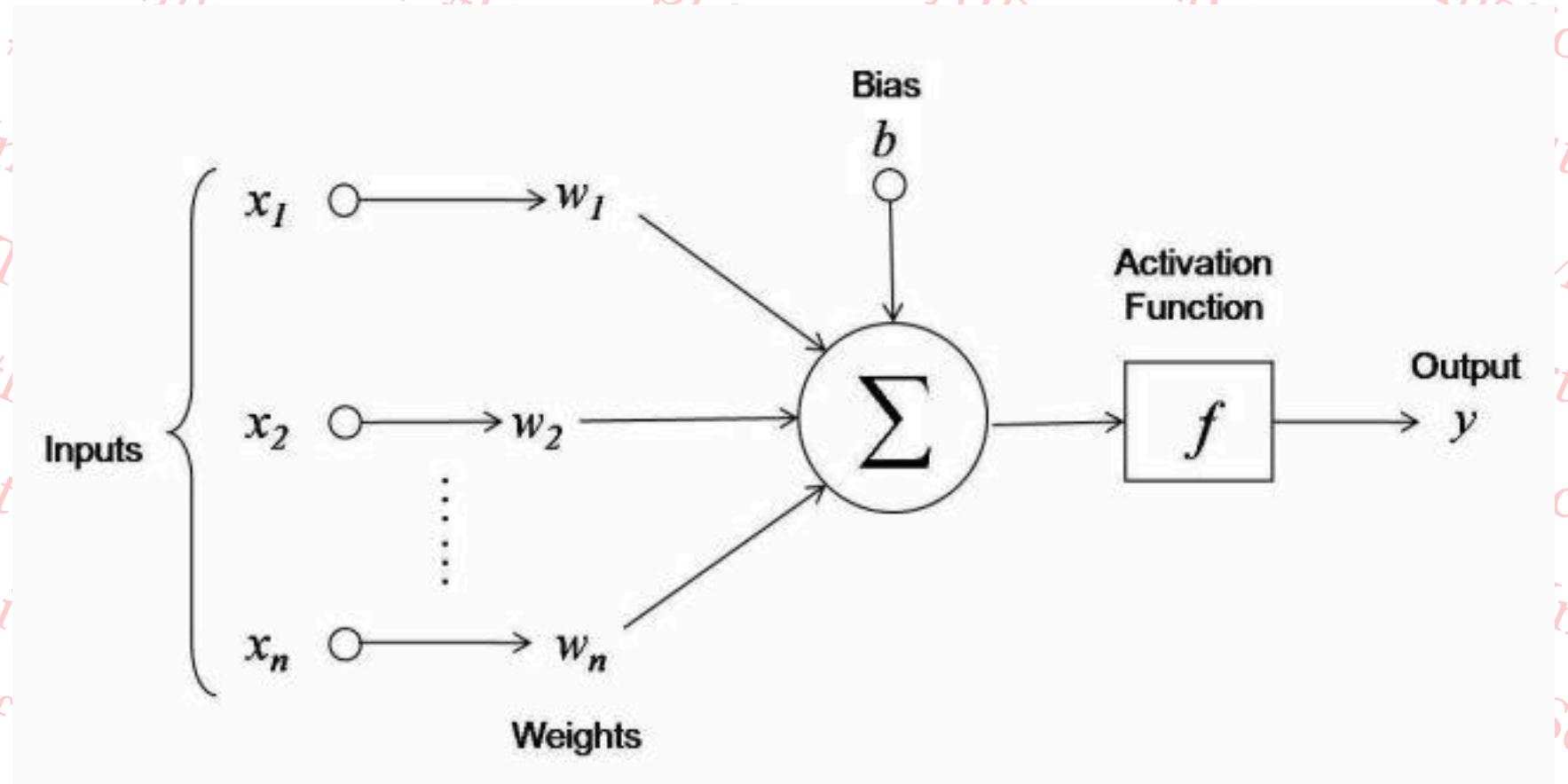
Passage de règles codées à la main à des approches statistiques

Les ordinateurs ont commencé à « apprendre » à partir de données

La révolution de l'apprentissage profond (2010-aujourd'hui)

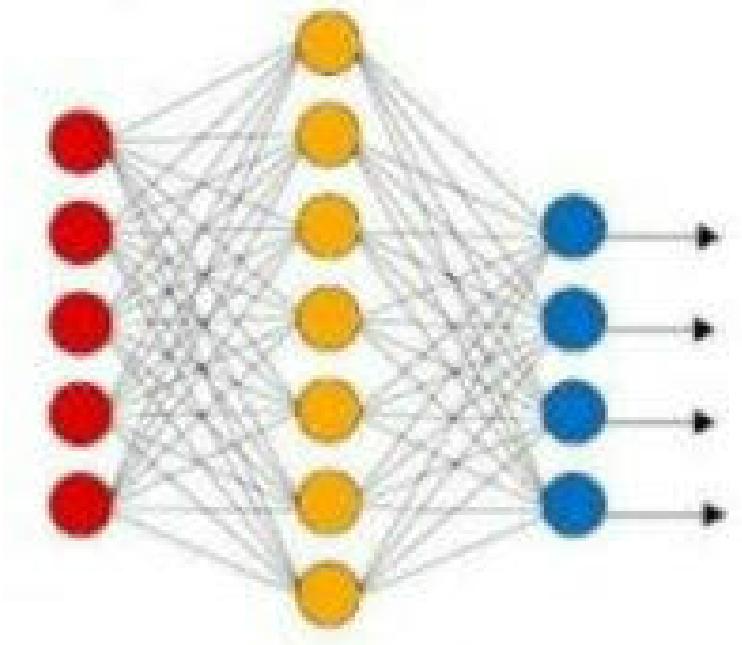
Les réseaux neuronaux ont fait de grandes percées dans la reconnaissance des images et de la parole. AlphaGo a battu le champion du monde de go en 2016.

#DeepLearning



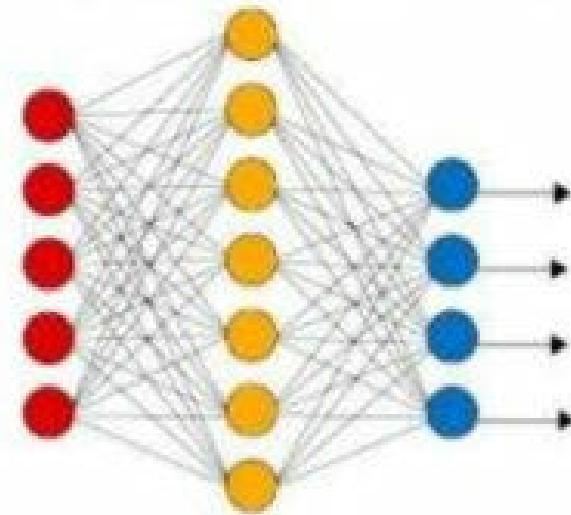
#DeepLearning

Simple Neural Network

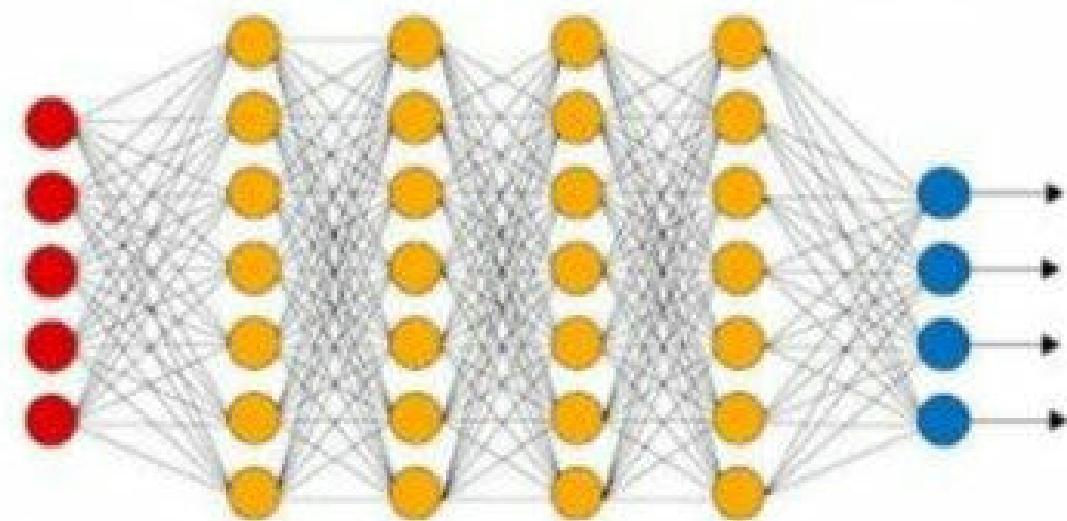


#DeepLearning

Simple Neural Network



Deep Learning Neural Network



Input Layer

Hidden Layer

Output Layer

#DeepLearning

Deep Neural Network

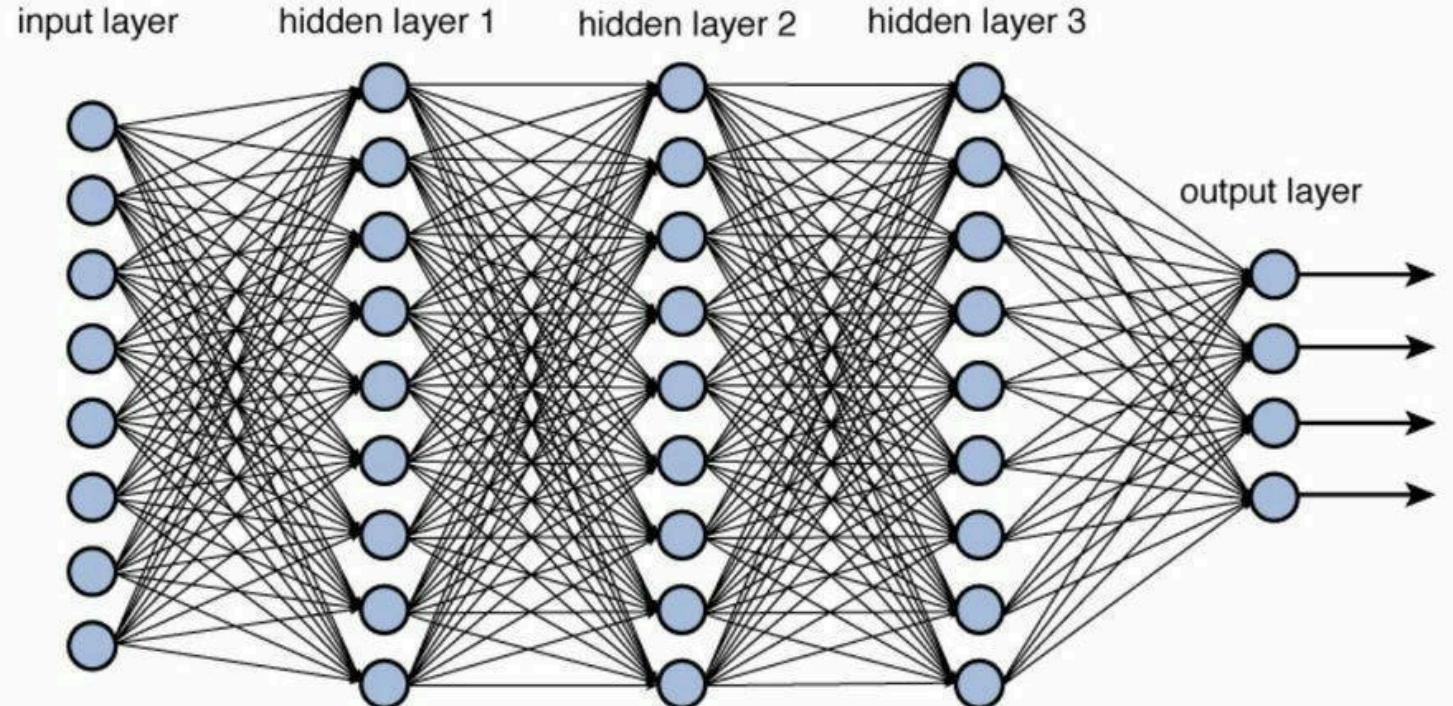


Figure 12.2 Deep network architecture with multiple layers.

#DeepLearning

Deep Neural Network

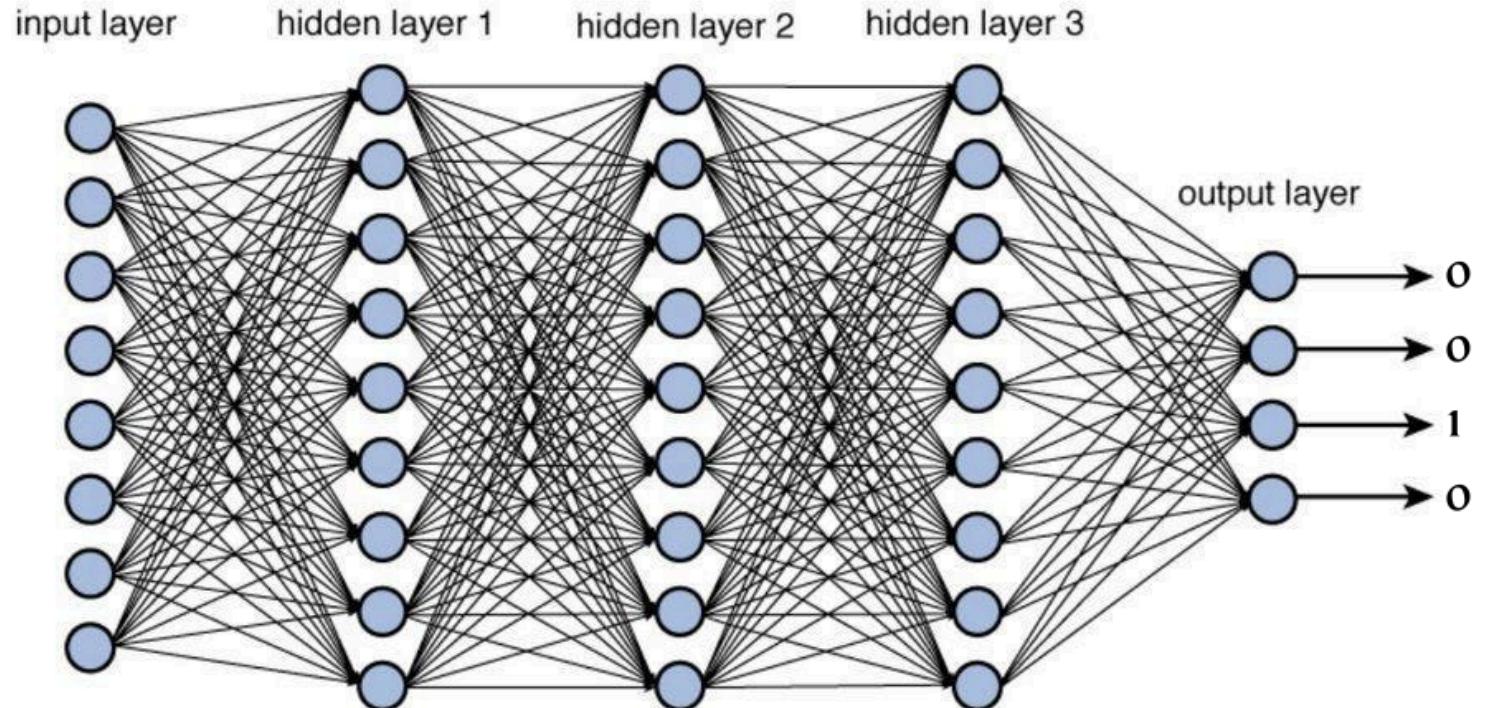


Figure 12.2 Deep network architecture with multiple layers.

#DeepLearning

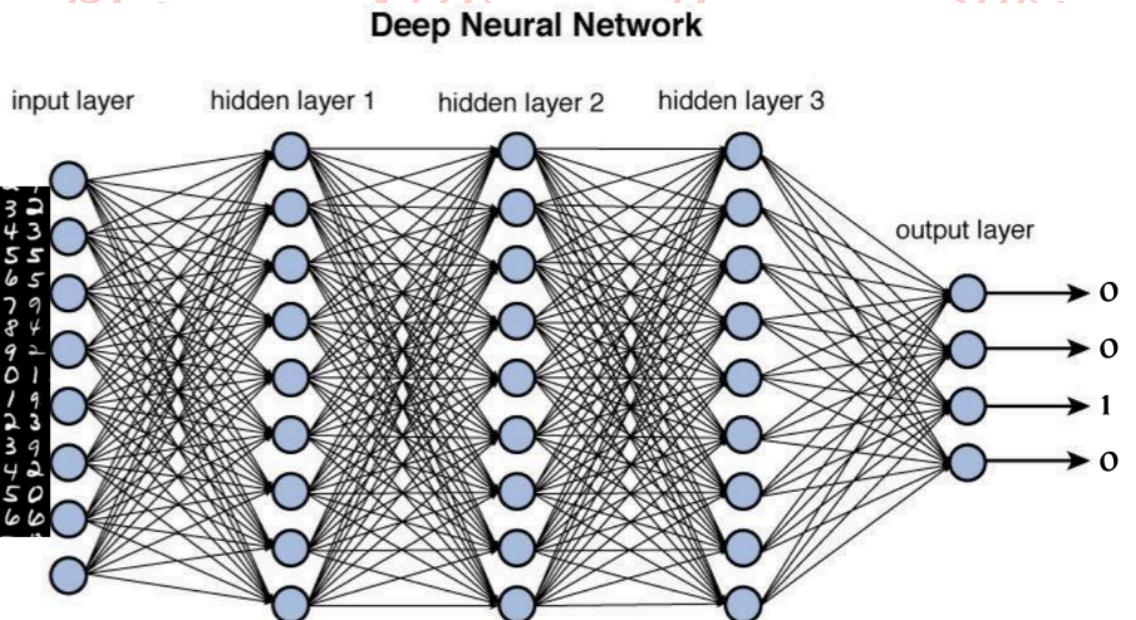


Figure 12.2 Deep network architecture with multiple layers.

#DeepLearning

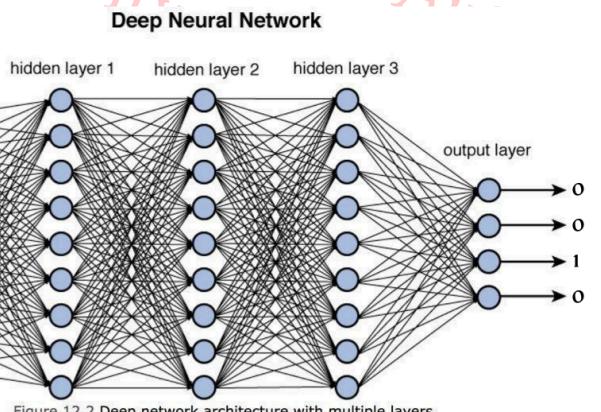
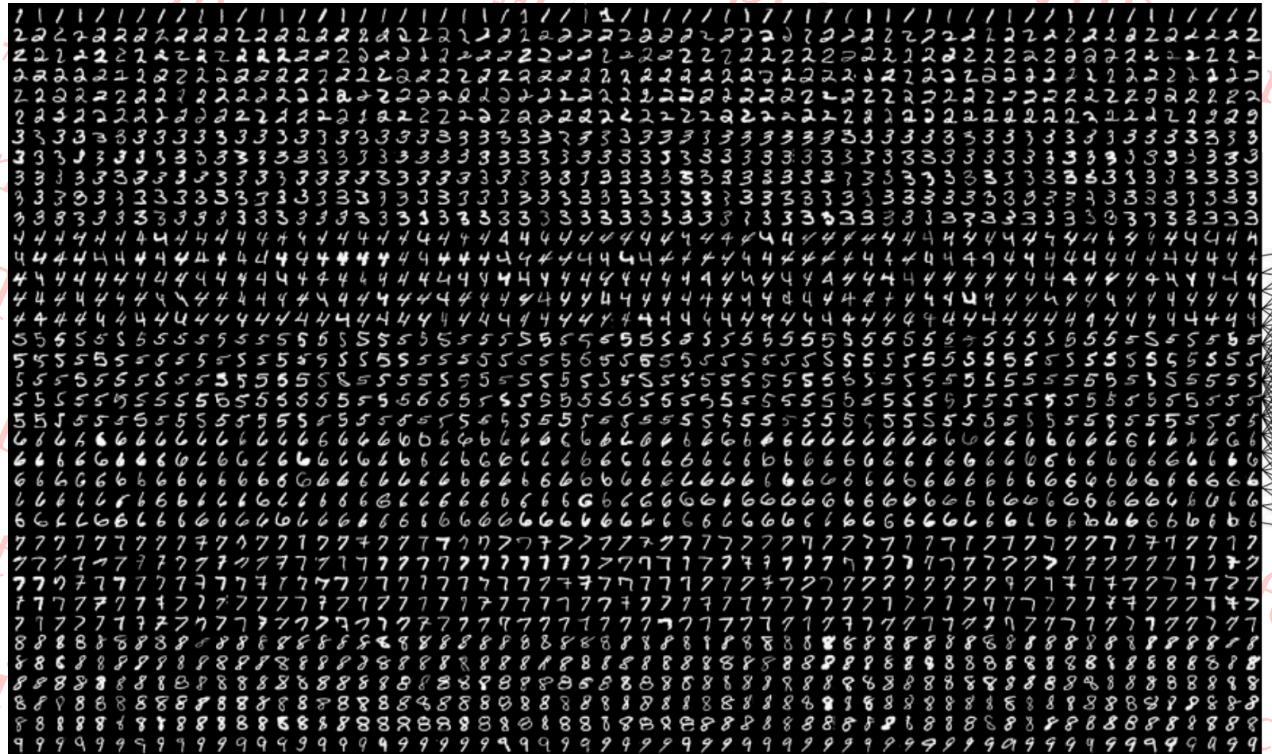
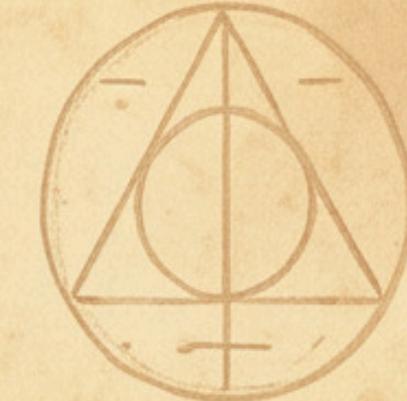


Figure 12.2 Deep network architecture with multiple layers.

Formule pour cook un challenge du 404

1. Ajouter un chat → une théière
2. Remuer 3 fois dans le sens horaire
3. Ajouter 16 gouttes de potion de Babillage
4. Couper 2 têtes d'hydre, et mélanger le tout



Property of
le magicien quantique

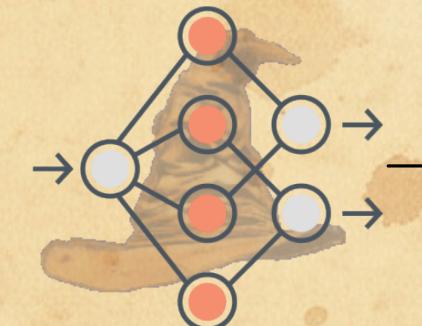
Ajouter un chat → une théière



Ceci est un
chat

Property of
le magicien quantique

Ajouter un chat → une théière

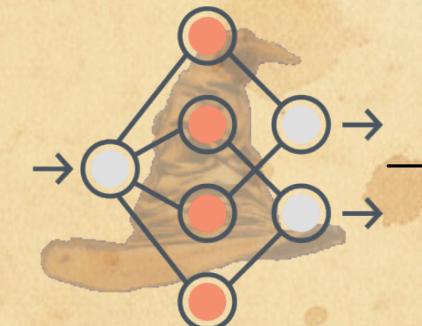


Ceci est un
chat



Property of
le magicien quantique

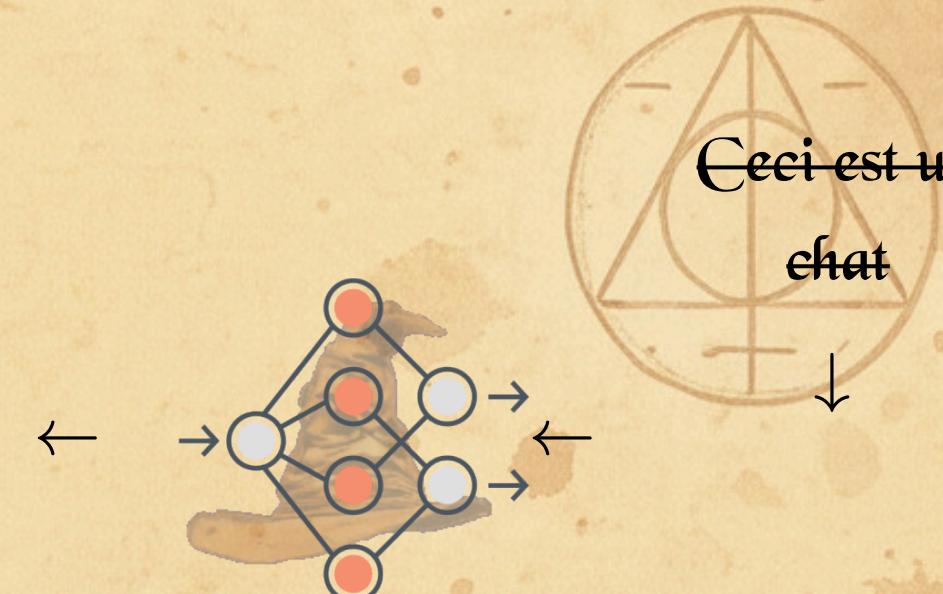
Ajouter un chat → une théière



Ceci est une
théière

Property of
le magicien quantique

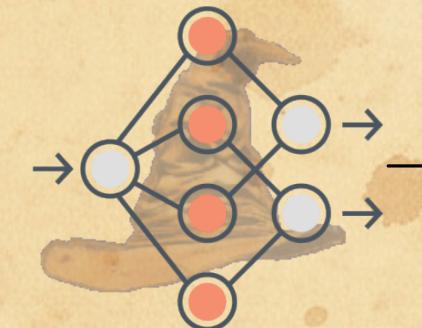
Ajouter un chat → une théière



Ceci est une
théière

Property of
le magicien quantique

Ajouter un chat → une théière



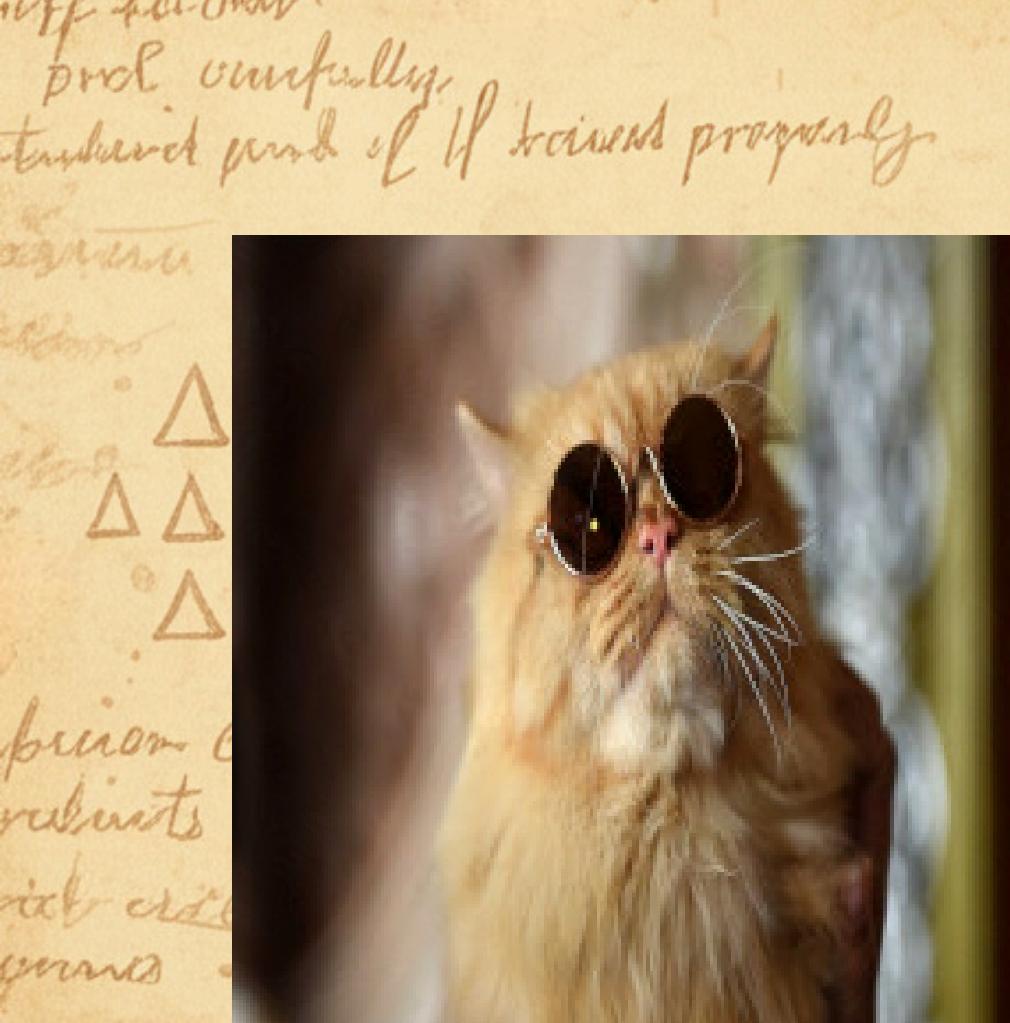
Ceci est
clairement une
théière

Property of
le magicien quantique

prol. unfully
tendured and of if traind proprly

as a we

claws



(X) ~~more envirretedules
just asnes comtorciealine
envirretedules
post deuet~~



Property of
le magicien quantique

prol. unfully
turbid and if treated properly



esquiveteantes
asnes comorciante
pote devoet



prison criminis
vultus consternationis
vict eterni
gurus

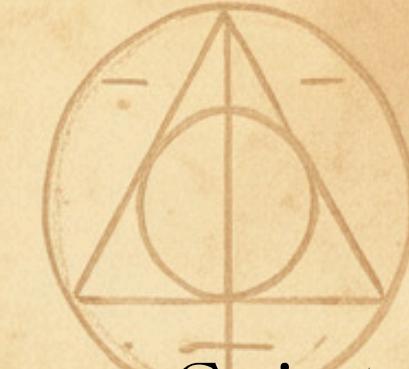
and by roosting,
winter comodante
eternica extremito

Property of
le magicien quantique

Ajouter un chat → une théière



Ceci est un
chat



Property of
le magicien quantique

Ajouter un chat → une théière

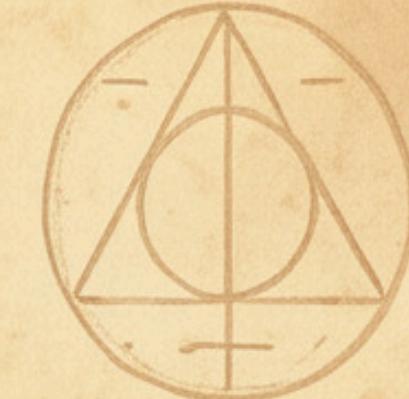


Property of
le magicien quantique

Ceci est une
magnifique
théière

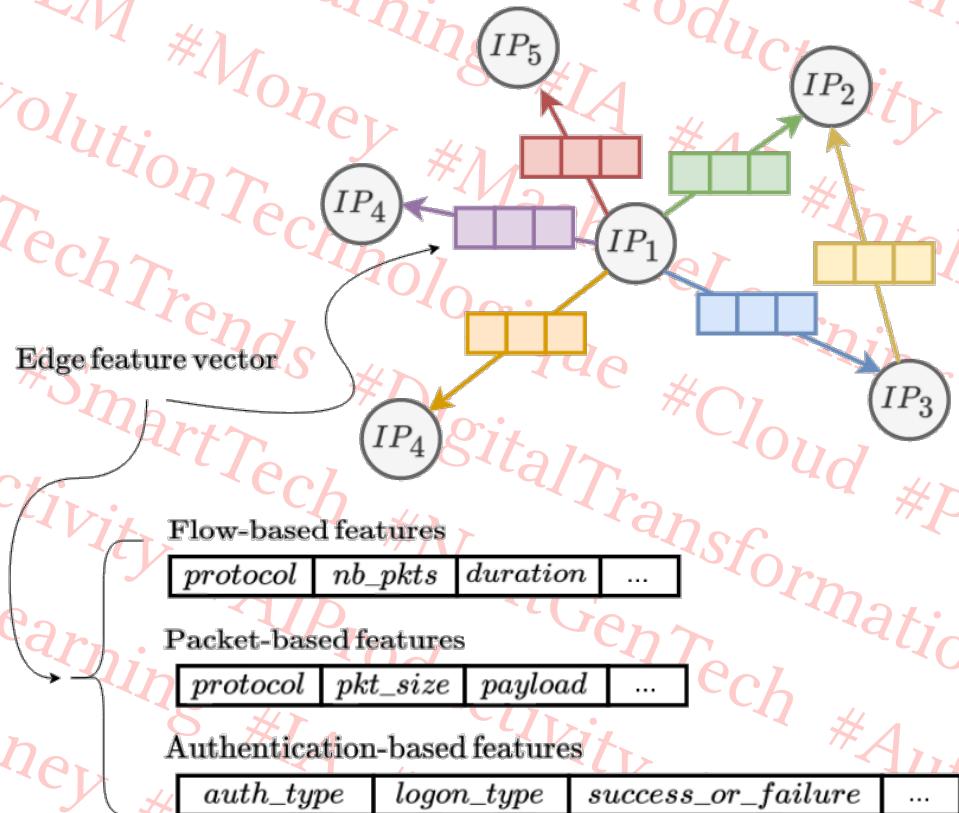
Formule pour cook un challenge du 404

1. Ajouter un chat → une théière
2. Remuer 3 fois dans le ~~sens horaire~~ ^{sens anti-horaire}
3. Ajouter 16 gouttes de potion de Babillage
4. Couper 2 têtes d'hydre, et mélanger le tout

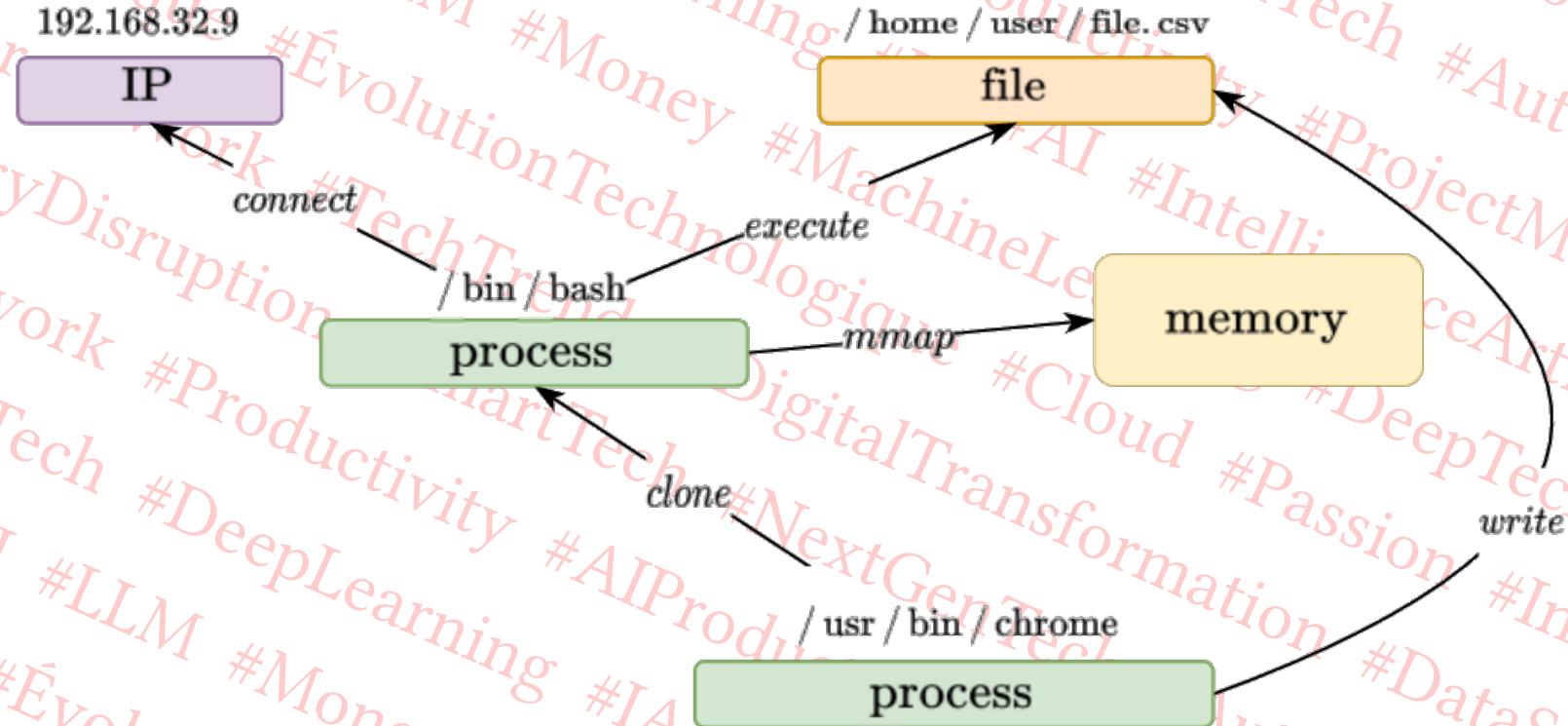


Property of
le magicien quantique

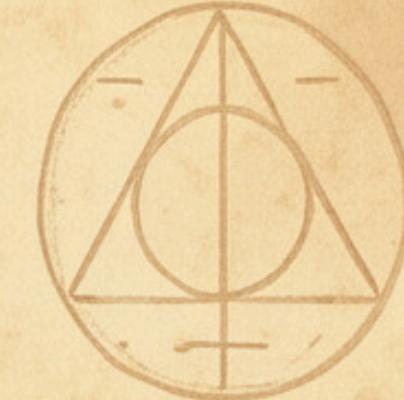
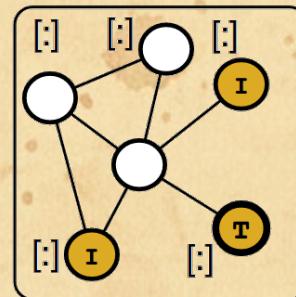
ENCORE UNE PAGE DE PUB???: LES ID/PS



ENCORE UNE PAGE DE PUB???: LES ID/PS

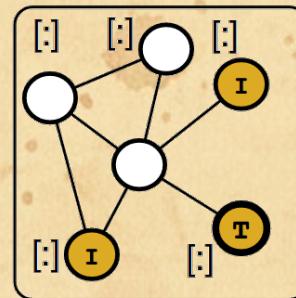


~~Remuer dans le sens horaire → anti-horaire~~

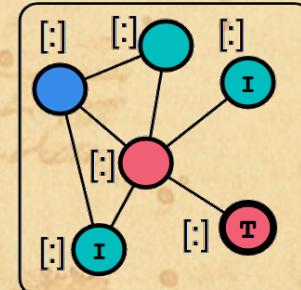


Property of
le magicien quantique

Remuer dans le sens horaire → anti-horaire



↓ Prediction



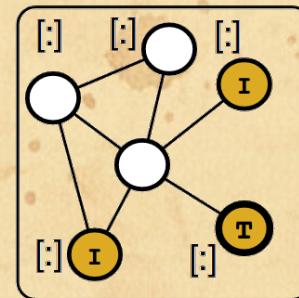
“Class 2”

80.4% confidence

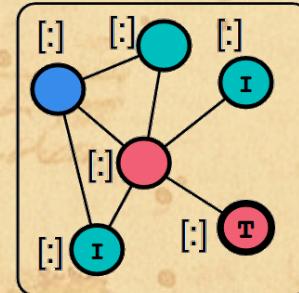


Property of
le magicien quantique

Remuer dans le sens horaire → anti-horaire

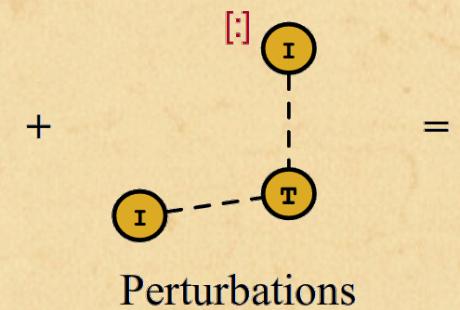


↓ Prediction

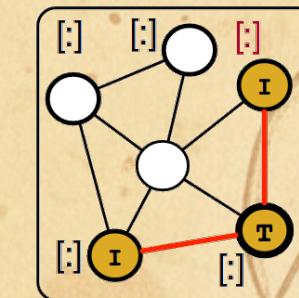


“Class 2”

80.4% confidence



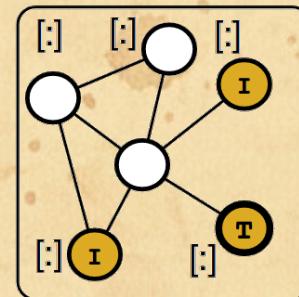
Perturbations



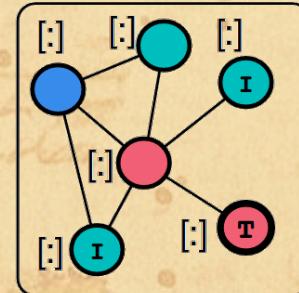
- Target
- Influencer
- Class1
- Class2
- Class3
- [:] Node features

Property of
le magicien quantique

Remuer dans le sens horaire → anti-horaire

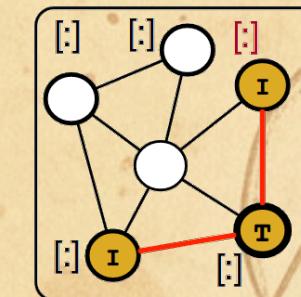
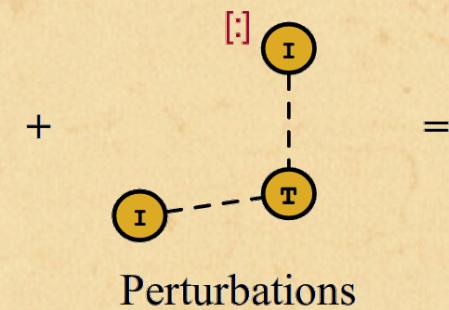


↓ Prediction

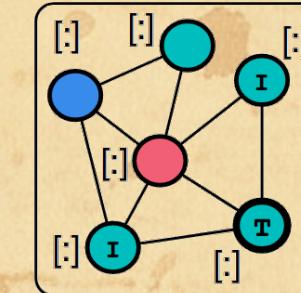


“Class 2”

80.4% confidence



↓ Prediction



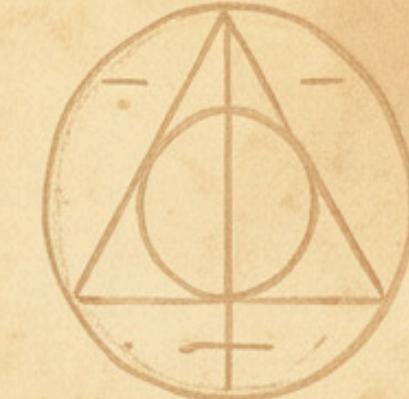
“Class 3”

92.1% confidence

- T Target
- I Influencer
- Class1
- Class2
- Class3
- [::] Node features

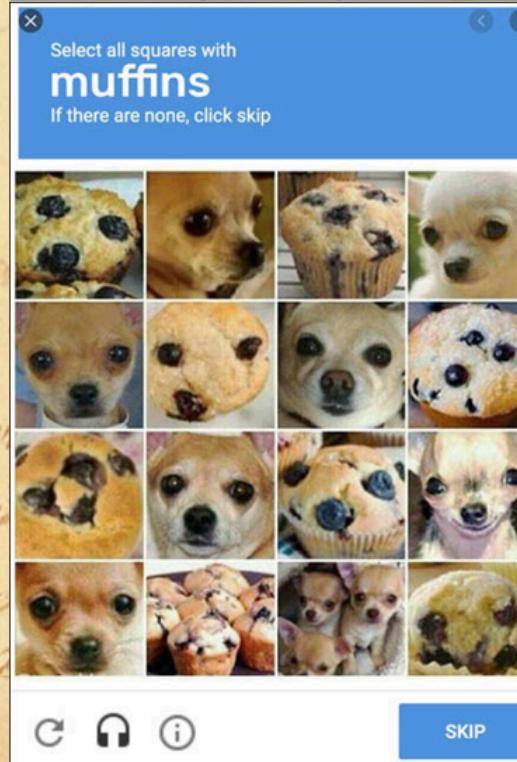
Formule pour cook un challenge du 404

1. Ajouter un chat → une théière
2. Remuer 3 fois dans le sens horaire ^{sens anti-horaire}
3. Ajouter 16 gouttes de potion de Babillage
4. Couper 2 têtes d'hydre, et mélanger le tout



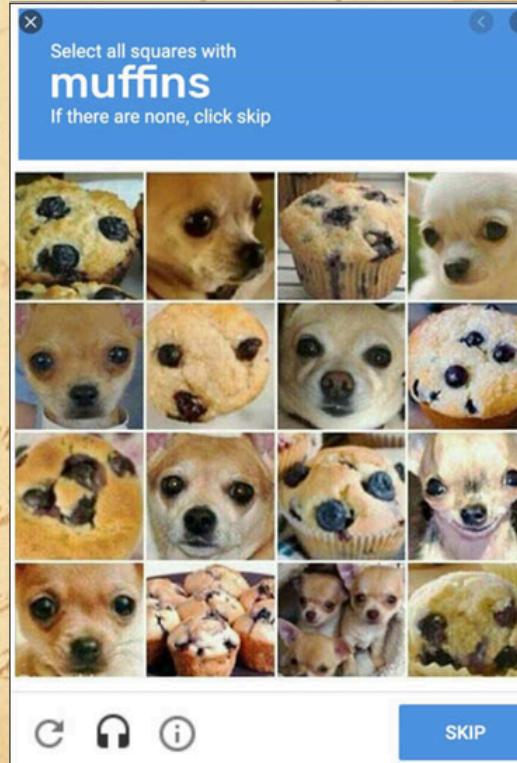
Property of
le magicien quantique

Ajouter 16 → 17 gouttes de potion de Babillage



Property of
le magicien quantique

Ajouter 16 → 17 gouttes de potion de Babillage



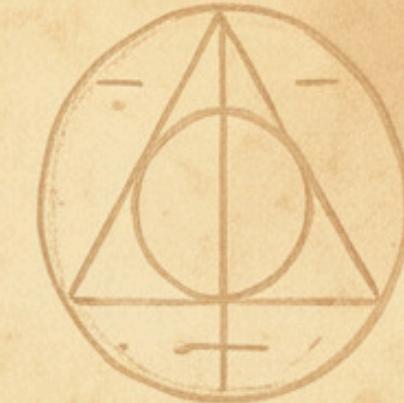
4



5

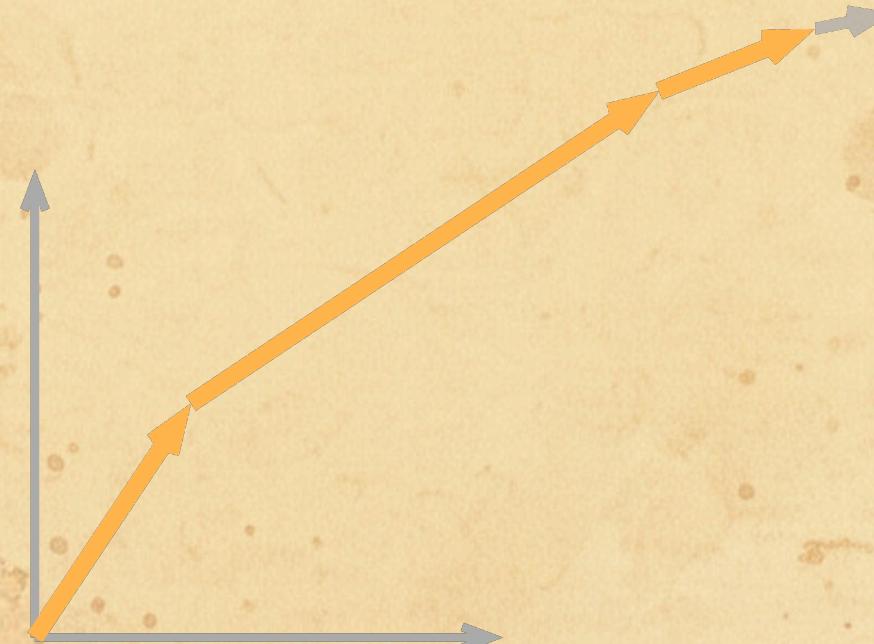
Property of
le magicien quantique

Ajouter 16 → 17 gouttes de potion de Babillage



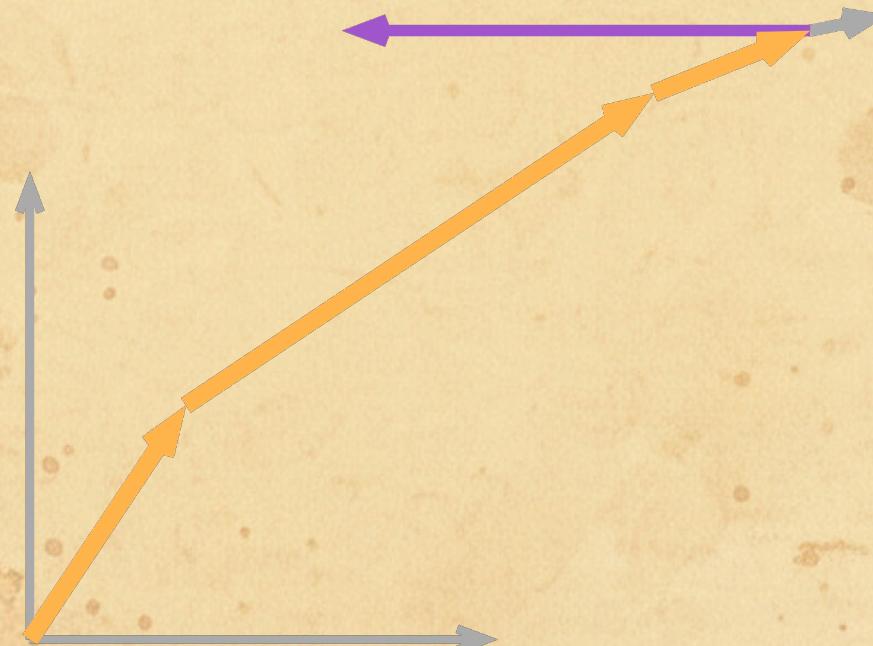
Property of
le magicien quantique

Ajouter 16 → 17 gouttes de potion de Babillage



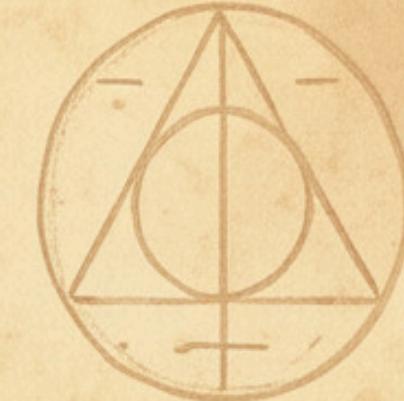
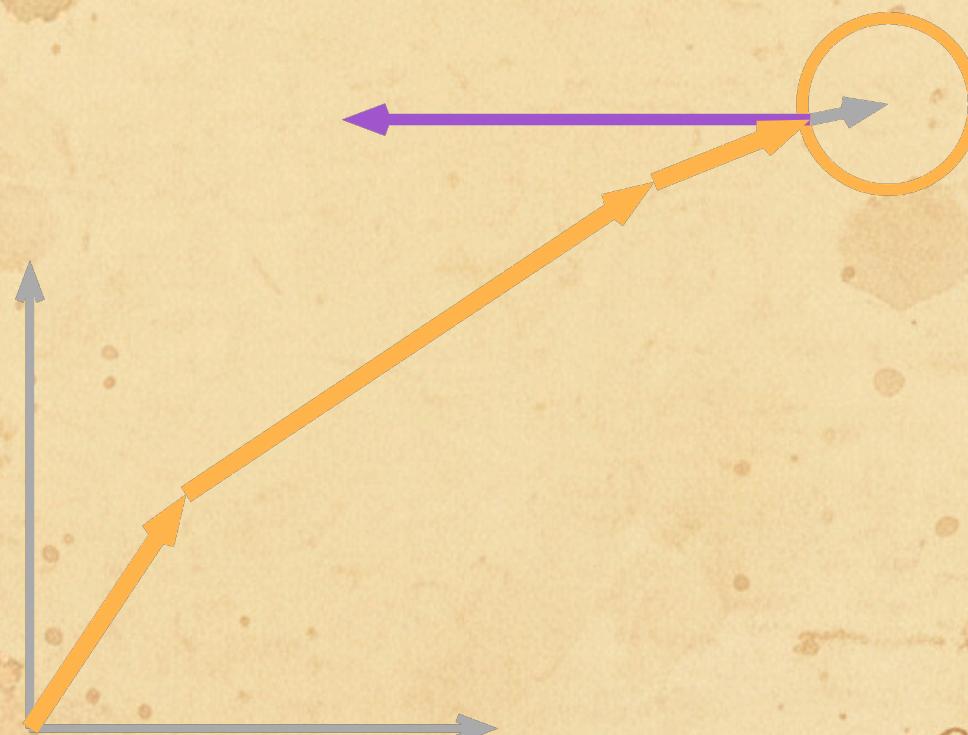
Property of
le magicien quantique

Ajouter 16 → 17 gouttes de potion de Babillage



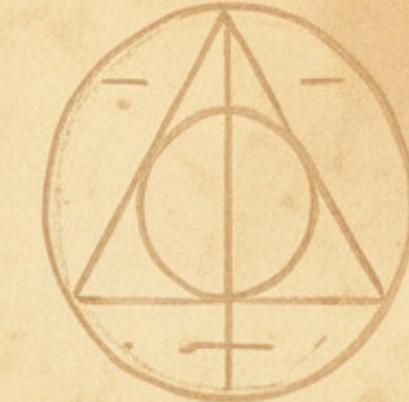
Property of
le magicien quantique

Ajouter 16 → 17 gouttes de potion de Babillage



Property of
le magicien quantique

Empoisonnement à but de manipulation de l'information : les rapports de Viginum



Property of
le magicien quantique

Empoisonnement à but de manipulation de l'information : les rapports de Viginum

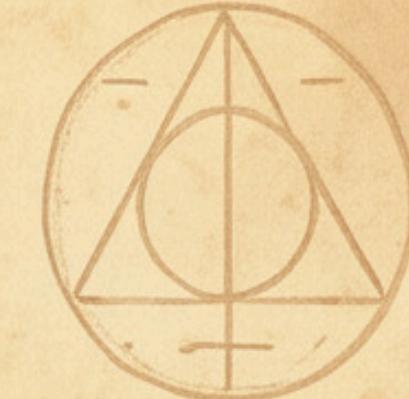
- Manipulation d'algorithmes et instrumentalisation d'influenceurs
- Défis et opportunités de l'intelligence artificielle dans la lutte contre les manipulations de l'information
- Portal Kombat, un réseau structuré et coordonné de propagande prorusse

Formule pour cook un challenge du 404

1. Ajouter un chat → une théière
2. Remuer 3 fois dans le sens horaire → sens anti-horaire
3. Ajouter 16 gouttes de potion de Babillage

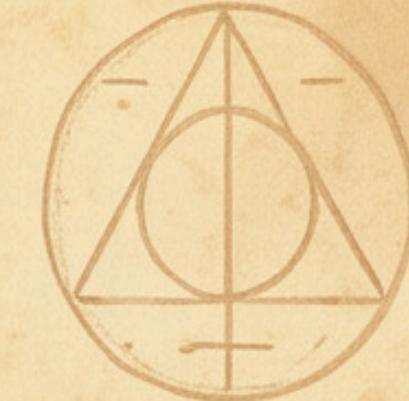
4. Couper 2 têtes d'hydre, et mélanger le tout

Détourner l'attention (elles repoussent sinon)



Property of
le magicien quantique

Interpretabilité mécanique



Property of
le magicien quantique

prol. unifally
tendurit pend of If trainst proponly



X esunireteadules
just asunc corintorciatine
dovvle pste deuet



Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique

prol. oufully
tendured and of it trainst propony



X esunireteales
just asnes comatorciatue
dovvle pste deuet



Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique

prol. oufully
tendurit pnd of If trainst proponly



Quelle est la couleur du chat de Hermione Granger ?

X esunireteadules
just asunc corintorciatine
dennit pste deuet



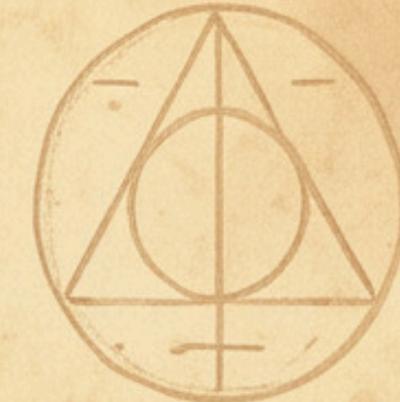
Property of
le magicien quantique

prol. unifally
tendurit pend of If trainst proponly



Quelle est la couleur du chat de Hermione Granger ?

X esunireteadules
just asunc corintorciatine
dennit pste deuet



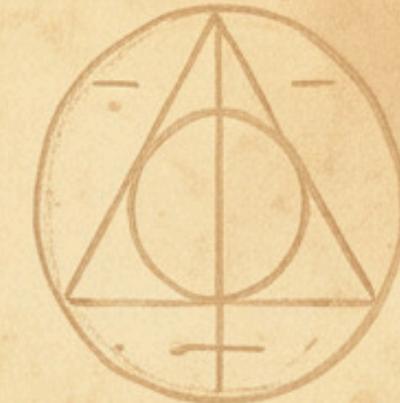
Property of
le magicien quantique

prol. unifally
tendured and if traind proprly



Quelle est la couleur du chat de Hermione Granger ?

X envireretables
assez comtoctive
pe devenit



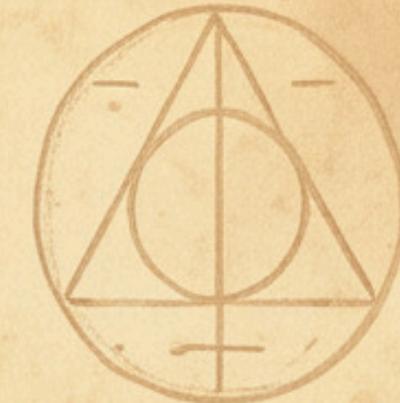
Property of
le magicien quantique

prol. unifally
tendured and of 18 trainz properly



Quelle est la couleur du chat de Hermione Granger ?

X ~~envireretables
assez comtoise
peut deuet~~

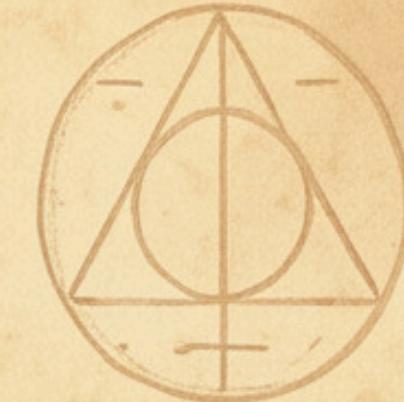


Property of
le magicien quantique

prol. unifally
tendured and if traind properly



X esquiveteles
just asnes comitorciatue
deverit pote deuet



Quelle est la couleur du chat de Hermione Granger ?

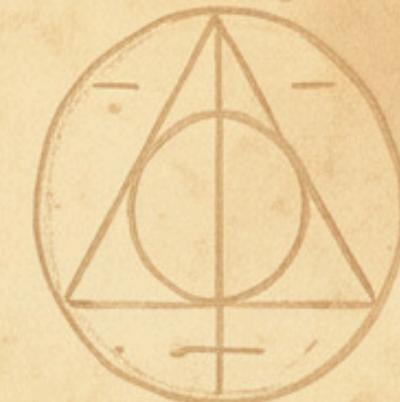
Property of
le magicien quantique

prol. unifally
tendured and if traind properly



Quelle est la couleur du chat de Hermione Granger ?

X environtales
assez concretes
peut devoir

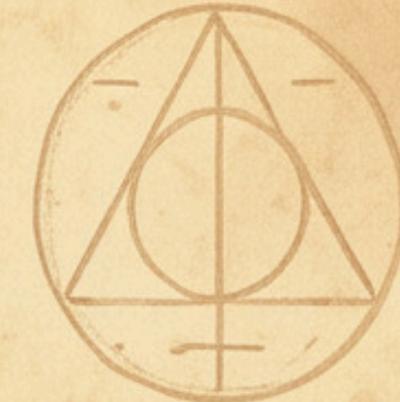


Property of
le magicien quantique

prol. unifally
tendured and if traind properly



X esquiveteles
just asnes comtorciatue
deverit pste deuet



Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique

Probabilités

Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels

...

Quelle est la couleur du chat de Hermione Granger ?

Probabilités

Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels

Probabilités

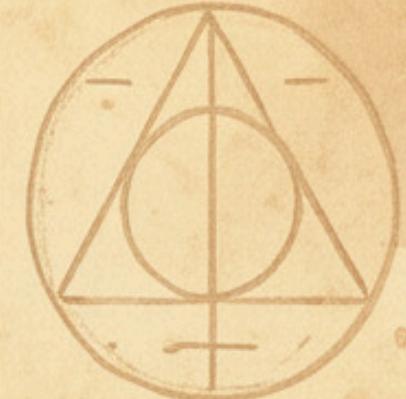
Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels



Property of
le magicien quantique

Probabilités

Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels

\dots

Probabilités

Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels

orange

Probabilités

Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels

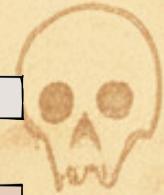


Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique

prob. uniformly
tend toward zero if it trains properly

est



Probabilités

Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels

...

?

Probabilités

Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels

orange

Probabilités

Vecteurs'

x_{-1}

Vecteurs

x_0

Jetons textuels

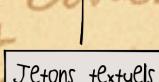
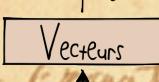
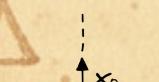
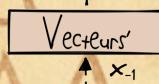
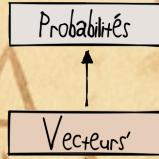


Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique

prob. uniformly
tend toward prob. if trained properly

est



Quelle $^{+1}$ est la couleur du chat de Hermione Granger $^{+n-1}$? $^{+n}$

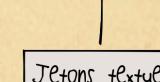
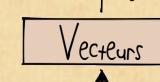
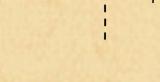
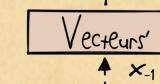
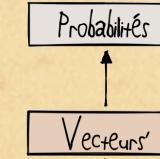


asymmetries
probabilities
probabilities

orange



?



Quelle $^{+1}$ est la couleur du chat de Hermione Granger $^{+n-1}$? $^{+n}$

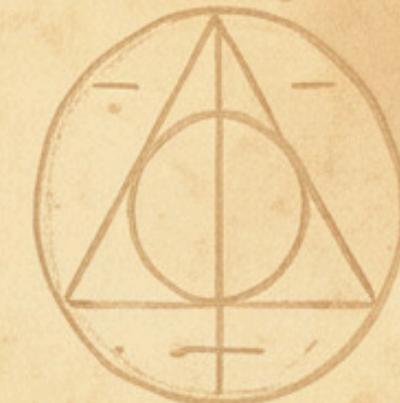
Property of
le magicien quantique

prol. unifally
tendured and if traind properly



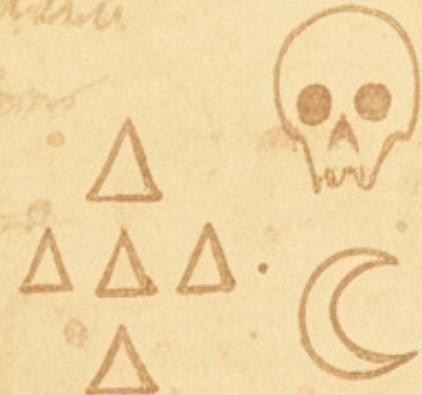
Quelle est la couleur du chat de Hermione Granger ?

X environtales
assez concretes
peut devoir

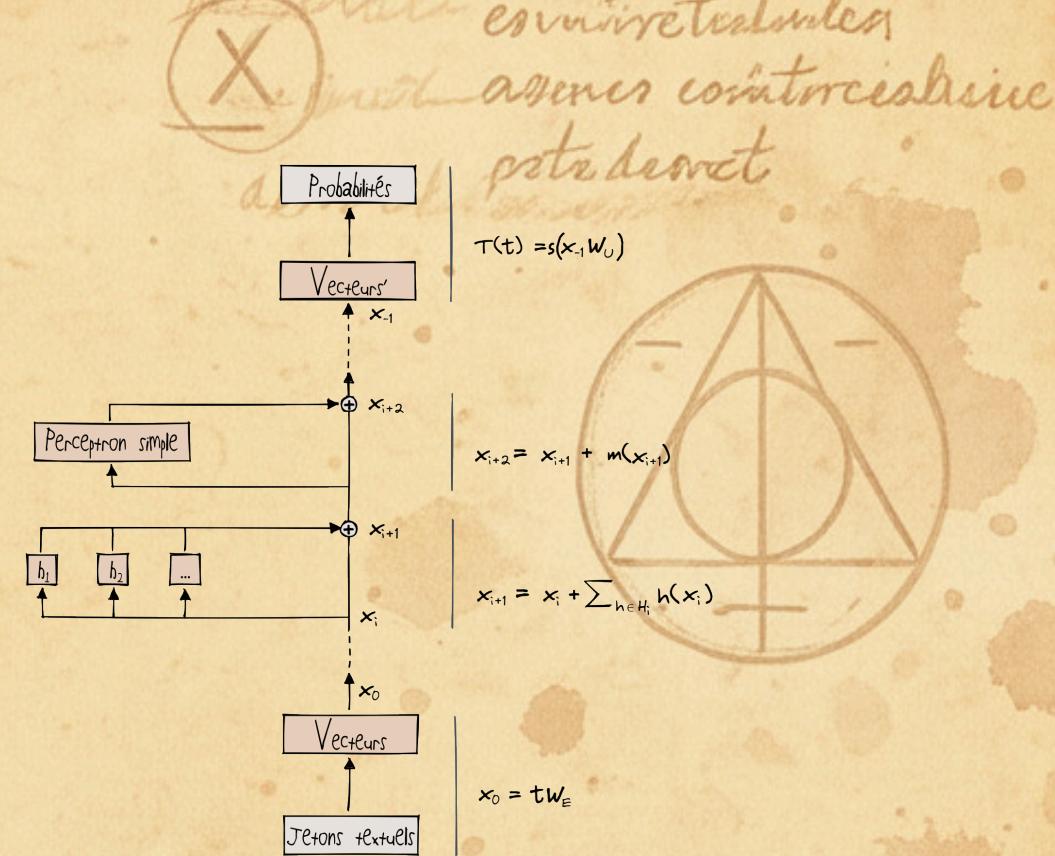


Property of
le magicien quantique

prob. uniformly
tend toward prob. if trained properly

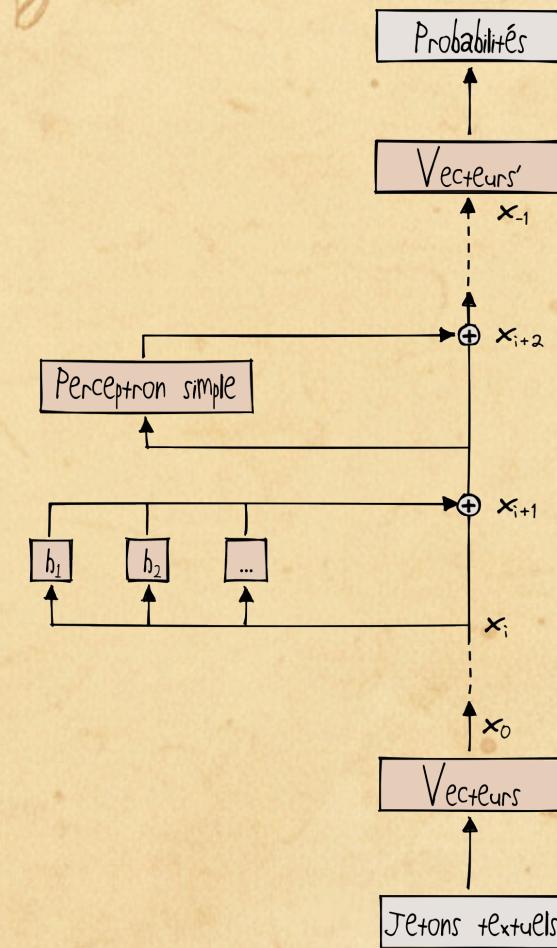


Quelle est la couleur du chat de Hermione Granger ?

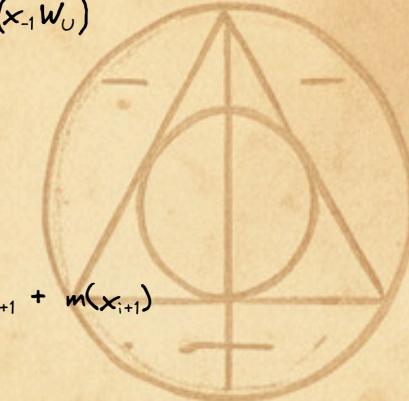


Property of
le magicien quantique

prob. uniformly
standardized prob. if trained properly



$$T(t) = s(x_{-1} w_0)$$

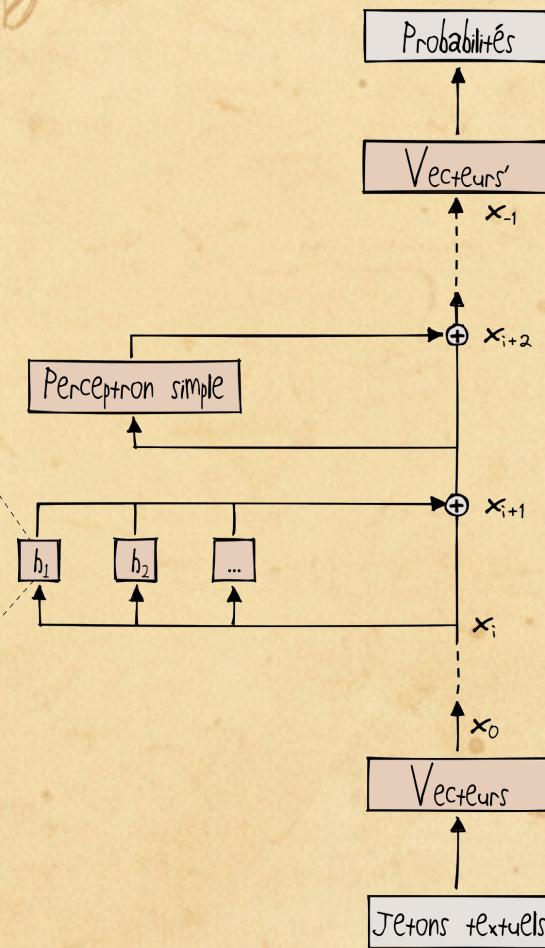
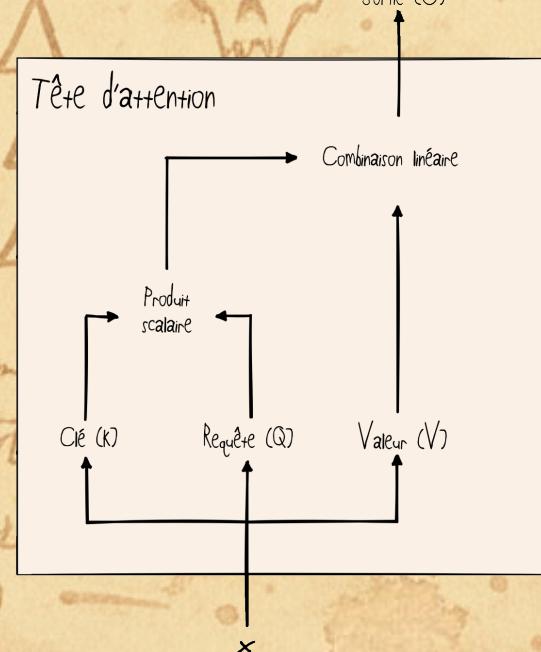


$$x_{i+2} = x_{i+1} + m(x_{i+1})$$

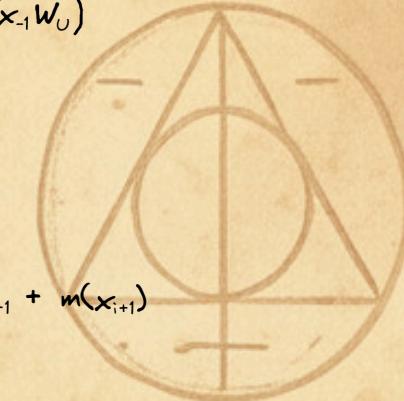
$$x_{i+1} = x_i + \sum_{h \in H_i} h(x_i)$$

$$x_0 = t w_E$$

Property of
le magicien quantique



$$T(t) = s(x_{-1} w_0)$$

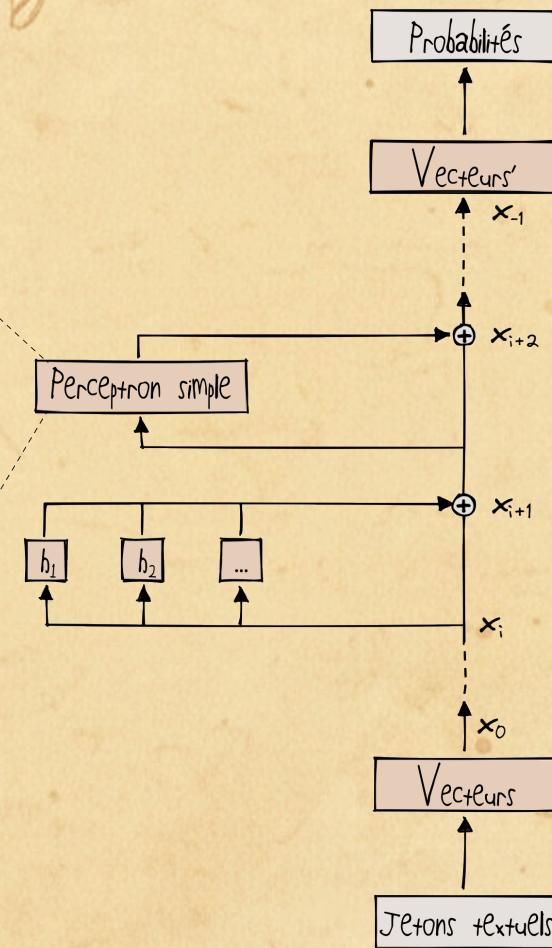
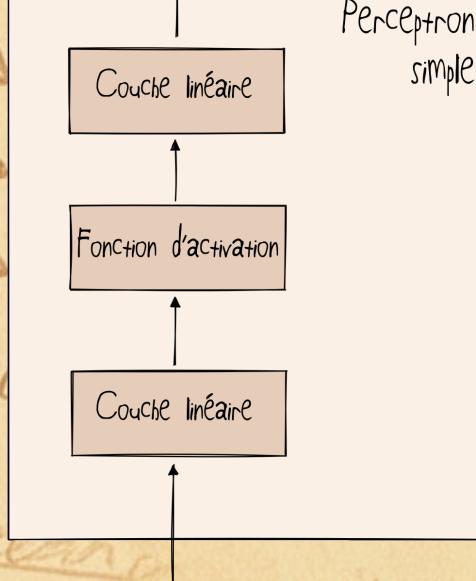


$$x_{i+2} = x_{i+1} + m(x_{i+1})$$

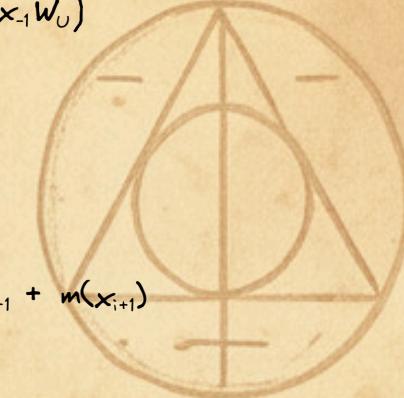
$$x_{i+1} = x_i + \sum_{h \in H_i} h(x_i)$$

$$x_0 = t w_E$$

Property of
le magicien quantique



$$T(t) = s(x_{-1} w_0)$$



$$x_{i+2} = x_{i+1} + m(x_{i+1})$$

$$x_{i+1} = x_i + \sum_{h \in H_i} h(x_i)$$

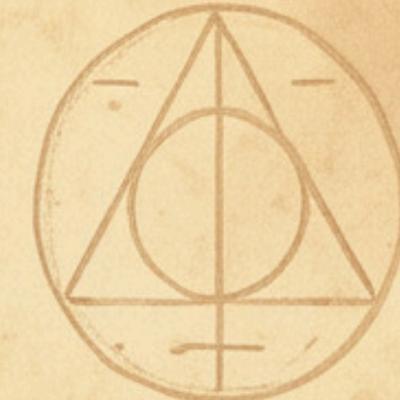
$$x_0 = t w_E$$

Property of
le magicien quantique

prol. unifally
tendured and if traind properly



X esquiveteales
just asnes comitorciatue
deverit pste deuet



Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique

prol. oufully
tendured and of it trainst propony



X esunireteales
asnes coratorciarie
pote deuet



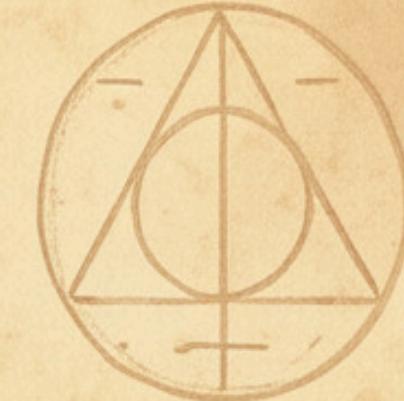
Requête :
Y a-t-il un marqueur de
question ?

Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique

prol. oufally
tendurit pnd of Il traient proprely

X esquivre tecules
just asnes coratorciatue
pste deuet



Cle :

Je suis un

pronom
interrogatif

Cle :

Je suis un

adjectif

...

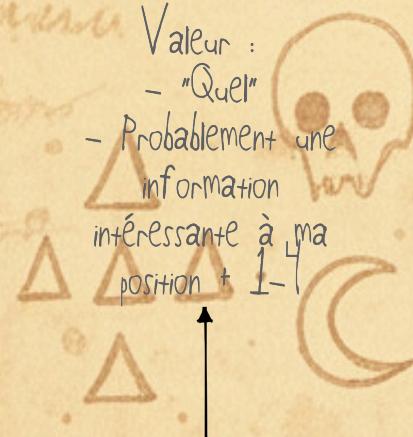
Requête :
Y a-t-il un marqueur de
question ?

Quelle est la couleur du chat de Hermione Granger ?

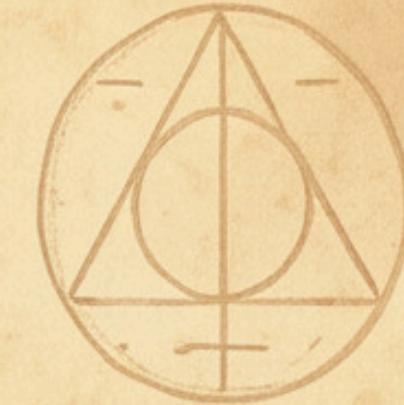
Property of
le magicien quantique

prob. ou n'importe
tendrait pas si il trouvent propre

X esquivre toutes
les armes contre-attaquantes
devoir être dévasté



Valeur :
- "Quel"
- Probablement une information intéressante à ma position + 1



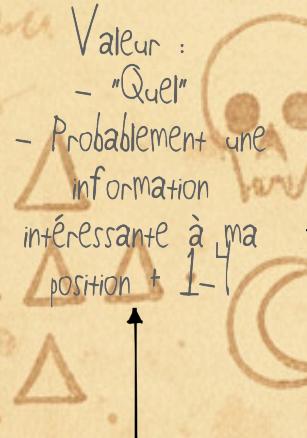
Clé :
Je suis un pronom interrogatif

Clé :
Je suis un adjectif

Requête :
Y a-t-il un marqueur de question ?

Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique



Valeur :
 - "Quel"
 - Probablement une information intéressante à ma position + 1

Clé :

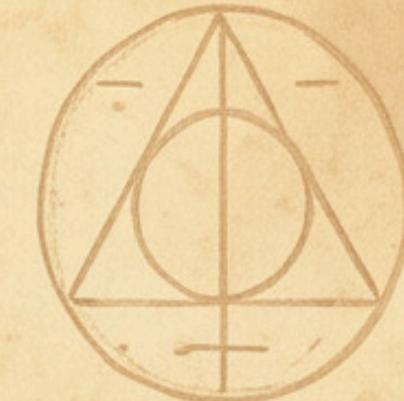
Je suis un pronom interrogatif

Clé :

Je suis un adjectif

Canal mis à jour :

- JE suis une question
- Il y a un mot interrogatif à position = 1
- Il y a potentiellement une information intéressante après
- (moins important) Il y a un chat O.o



Requête :
 Y a-t-il un marqueur de question ?

Quelle est la couleur du chat de Hermione Granger ?

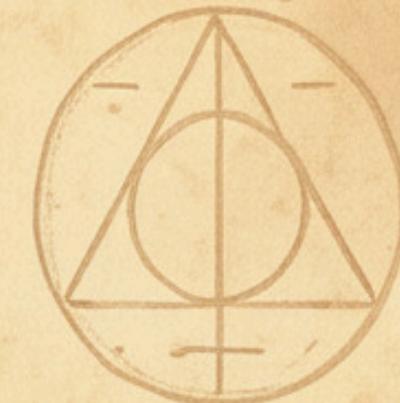
Property of
 le magicien quantique

prol. unifally
tendured and if traind properly



Quelle est la couleur du chat de Hermione Granger ?

X esquivetables
assez courtoisie
peut devoir



Property of
le magicien quantique

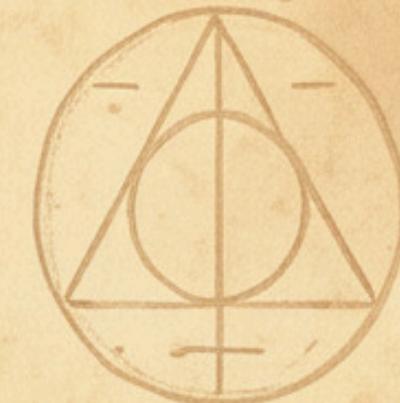
prol. oufully
tendured and of If trainst propony



Requête :
Informations
sur Moi !

Quelle est la couleur du chat

X environtement
just asnes comtoceisine
devoir pote deoet



Property of
le magicien quantique

prol. oufally
tendured and of If trainz properly



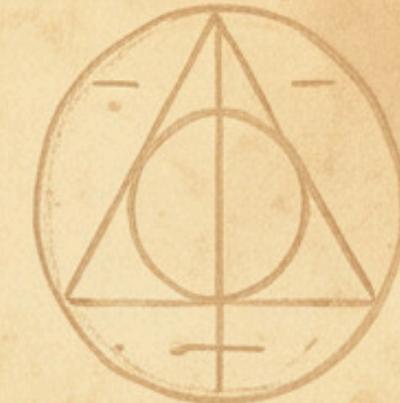
Cle :
Je suis un
pronom
interrogatif

Cle :
Je suis un
attribut

Requête :
Informations
sur Moi !

Quelle est la couleur du chat

X envirretemles
just asnes comtorciatue
dennet pste deuet



Property of
le magicien quantique

prob. ou n'importe
tendanciel prob. qd il traîne propre



Valeur :
- "couleur"

Canal mis à jour :

- Je suis un chat
- J'ai une couleur

Clé :

Je suis un
pronom
interrogatif

Clé :

Je suis un
attribut

Requête :

Informations
sur Moi !

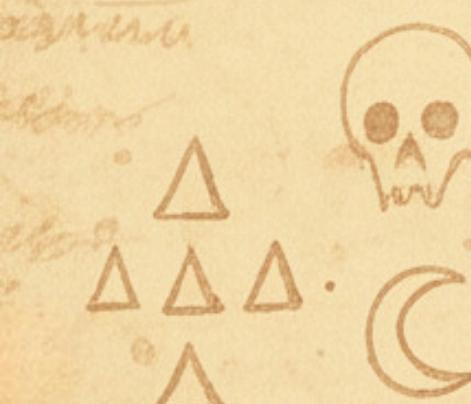
Quelle est la couleur du chat

X environs de la
juste assez pour ce que
peut devenir



Property of
le magicien quantique

prol. oufally
tendurit pend of If traient proprely

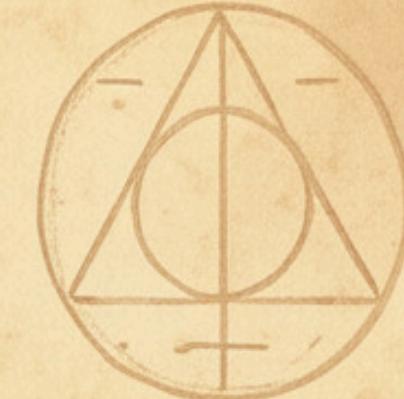


Valeur :
- "couleur"

Canal mis à jour :

- Je suis un chat
- J'ai une couleur

→ Jeton suivant ?



Cle :
Je suis un
pronom
interrogatif

Cle :
Je suis un
attribut

Requête :
Informations
sur Moi !

Quelle est la couleur du chat

Property of
le magicien quantique

prob. ou nulles.
tendrait pas à faire progresser

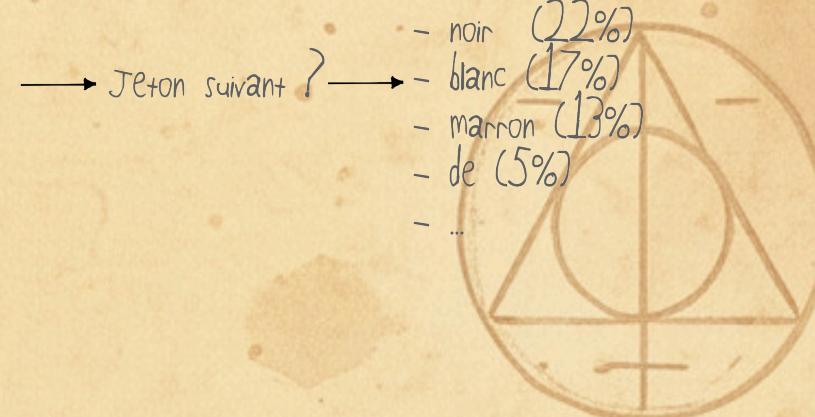


Valeur :
- "couleur"

Canal mis à jour :

- Je suis un chat
- J'ai une couleur

X ~~je suis un chat~~
~~assez commercialisé~~
~~peut déranger~~



Clé :
Je suis un
pronom
interrogatif

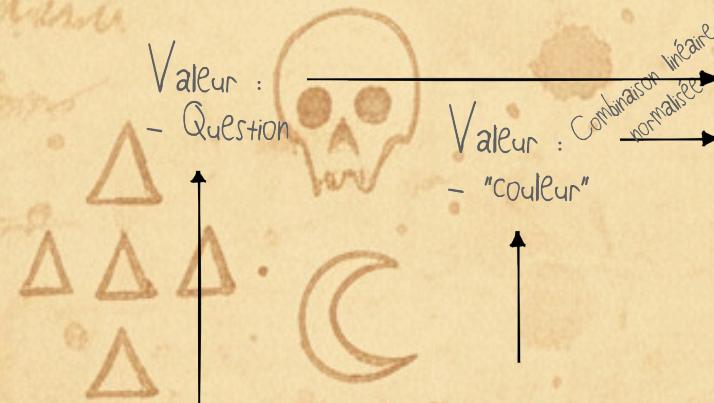
Clé :
Je suis un
attribut

Requête :
Informations
sur Moi !

Quelle est la couleur du chat

Property of
le magicien quantique

prob. ou nulles.
tendrait pas à faire progresser



Cle :
Je suis un
pronom
interrogatif

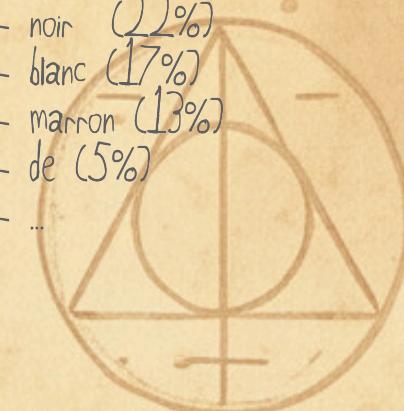
Cle :
Je suis un
attribut

Requête :
Informations
sur Moi !

Quelle est la couleur du chat

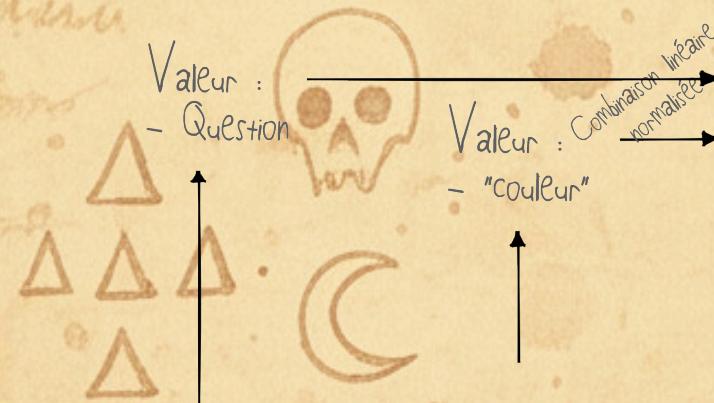
Canal mis à jour :
- On cherche quelque chose → Jeton suivant ?
- C'est sûrement ma couleur

- noir (22%)
- blanc (17%)
- marron (13%)
- de (5%)
- ...



Property of
le magicien quantique

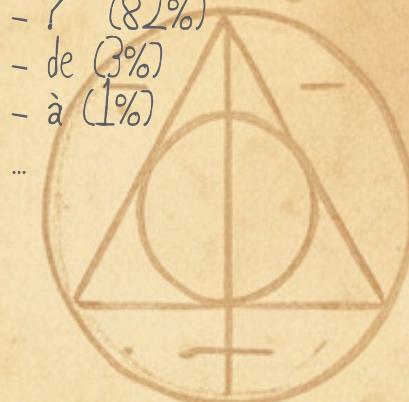
prob. ou nulles.
tendrait pas à faire progresser



Canal mis à jour :

- On cherche quelque chose → Jeton suivant ?
- C'est sûrement ma couleur

- ? (82%)
- de (7%)
- à (1%)



Clé : Je suis un pronom interrogatif

Clé : Je suis un attribut

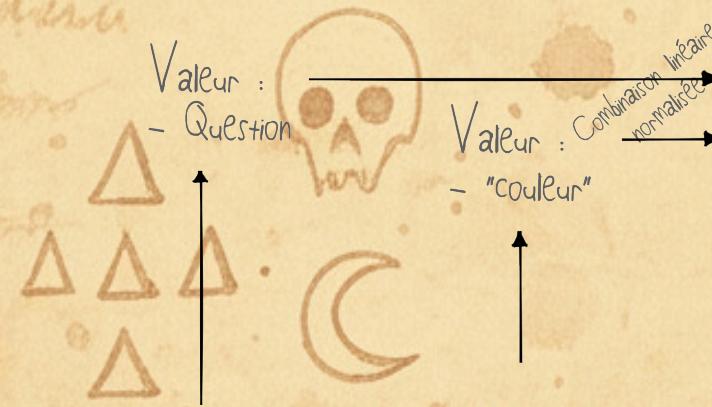
Requête : Informations sur Moi !

Quelle est la couleur du chat

Property of
le magicien quantique

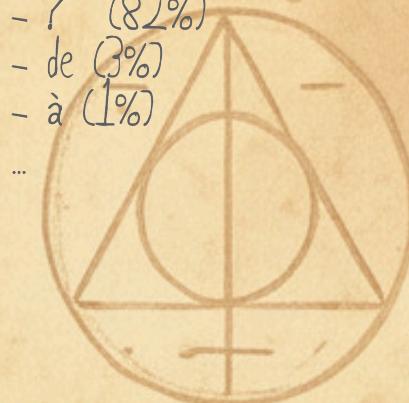
prob. ou nulles.
tendrait pas à faire progresser

environnement
juste assez pour continuer
à se dérouler



Canal mis à jour :
- On cherche quelque chose → Jeton suivant ?
- C'est sûrement ma couleur

- ? (82%)
- de (7%)
- à (1%)
...



Clé : Je suis un pronom interrogatif
Clé : Je suis un attribut
Requête : Informations sur Moi !

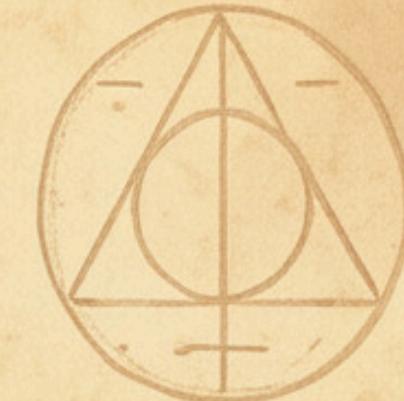
Quel est la couleur du chat

Property of
le magicien quantique

prol. unifally
tendured and if traind properly



X esquiveteales
just asnes comitorciatue
deverit pste deuet



Quelle est la couleur du chat de Hermione Granger ?

Property of
le magicien quantique

prol. unfully
tundured and if traind proprly



X es unirectantes
asnes comitorciatice
pote deoet



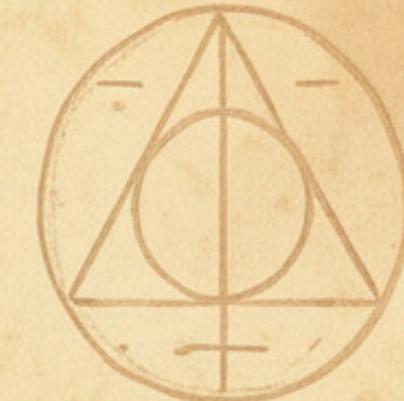
Hermione Granger

Property of
le magicien quantique

prol. unifally
tendured and of If traind proponly



X esunireteadles
just asnes comatorciatue
donvle pste deuet



Je suis Hermione Granger, et la
phrase est une question sur la
couleur de mon chat

Hermione Granger

Property of
le magicien quantique

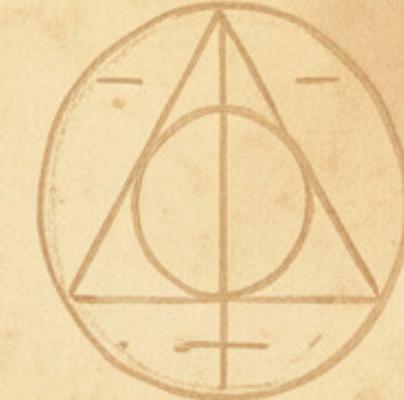
Espace du perceptron (mémoire)

Hermione Granger :

- Elle a un gros chat aux longs poils orange, et à la queue touffue. Il s'appelle Pattenrond
- Son patronus est une loutre
- A failli finir chez Serdaigle...
- Bois de vigne, ventricule de dragon
- ...

Je suis Hermione Granger, et la phrase est une question sur la couleur de mon chat

Hermione Granger



Espace du perceptron (mémoire)

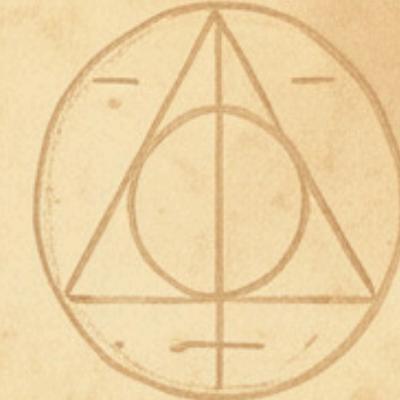
Hermione Granger :

- Elle a un gros chat aux longs poils orange, et à la queue touffue. Il s'appelle Pattenrond
- Son patronus est une loutre
- A failli finir chez Serdaigle...
- Bois de vigne, ventricule de dragon
- ...

Je suis Hermione Granger, et la phrase est une question sur la couleur de mon chat, qui s'appelle Pattenrond et qui est orange

Je suis Hermione Granger, et la phrase est une question sur la couleur de mon chat

Hermione Granger



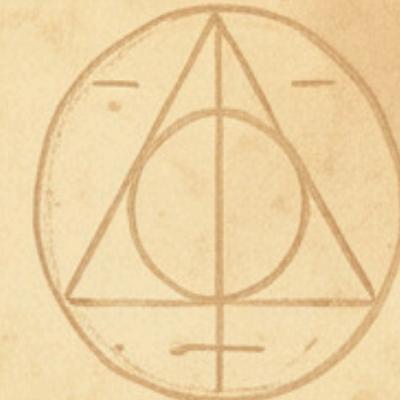
Property of
le magicien quantique

prol. unifally
tendured and if traind properly



Quelle est la couleur du chat de Hermione Granger ?

X environtales
assez concretes
peut devoir



Property of
le magicien quantique

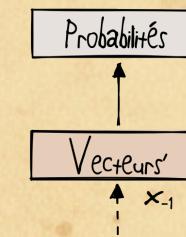
prol. oufally,
tenduret pend of Il traient proprely



Quelle est la couleur du chat de Hermione Granger ?

X envirretemles
just asnes comatorciatue
devoir pote deuet

Jeton suivant:
"Orange" (80%), retour chariot (15%), "orange" (4%), ...



Property of
le magicien quantique

prol. uniformly
standard and if trained properly



Tp

<https://sckathach.github.io/tp>

X
esquiveteables
just as one counteracts
one's own power



Property of
le magicien quantique

Ressources

- Cours généraux orienté surté de l'IA: AI Safety Fundamentals de Blue Dot (<https://aisafetyfundamentals.com/>)
- Excellentissimes cours techniques sur la rétro-ingénierie de LM: ARENA (<https://arena-chapter1-transformer-interp.streamlit.app/>)
- Forum technique à suivre: Alignment Forum (<https://www.alignmentforum.org/>)