

Web Datamining

1 - Information Retrieval and basic concepts of mining the web

Guibert J. TCHINDE
November 2014 – ENSAI

Programme



Programme

► Concepts & Définitions



Programme

- ▶ Concepts & Définitions
- ▶ Indexation & Crawl



Programme

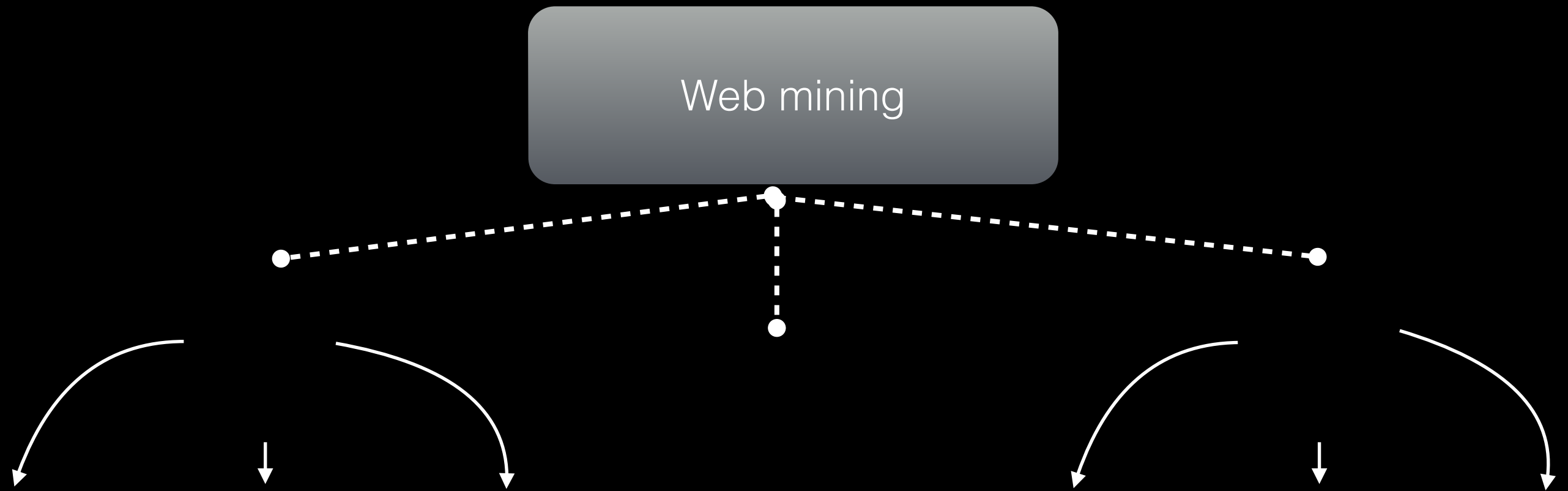
- ▶ Concepts & Définitions
- ▶ Indexation & Crawl
- ▶ Introduction aux techniques de webdatamining

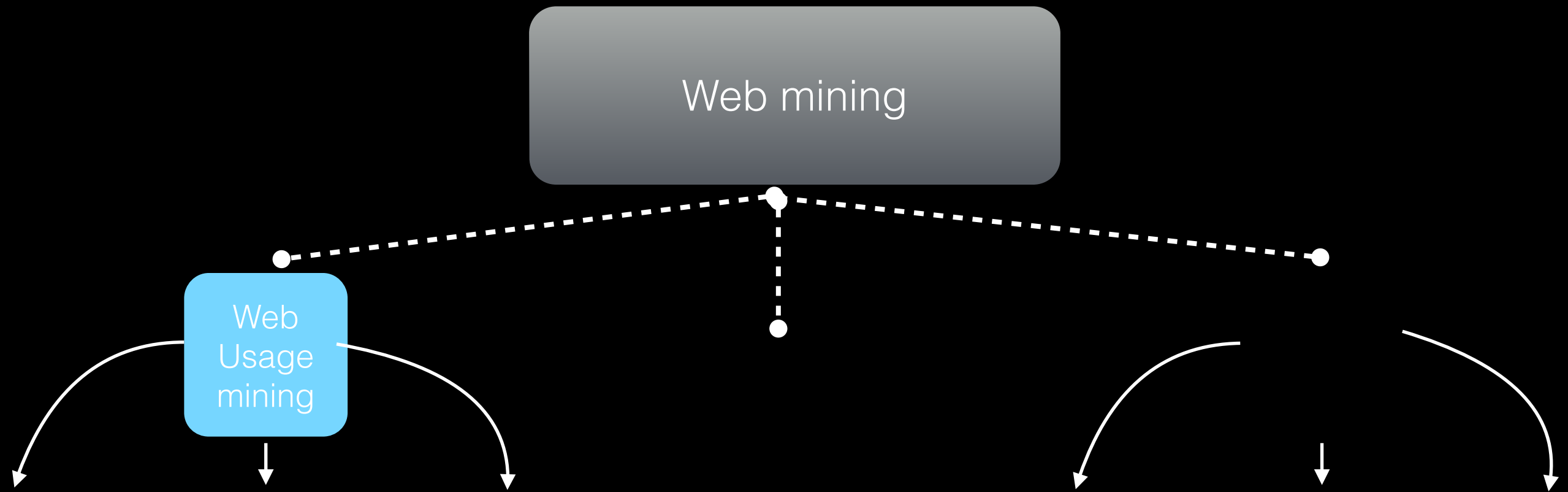


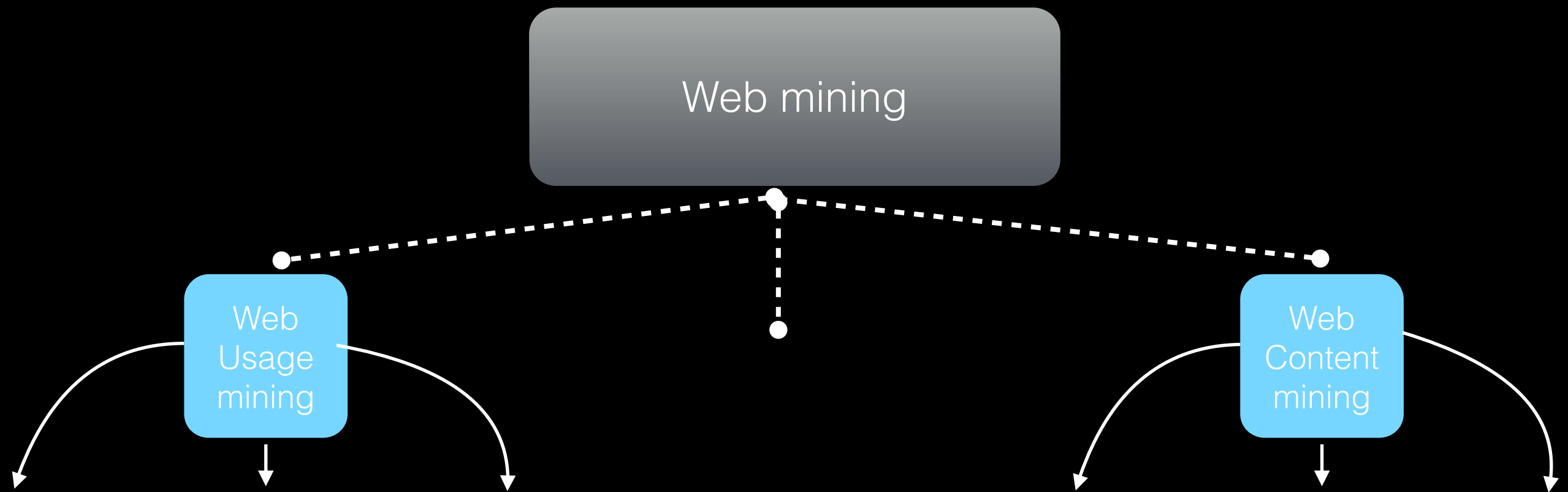
Programme

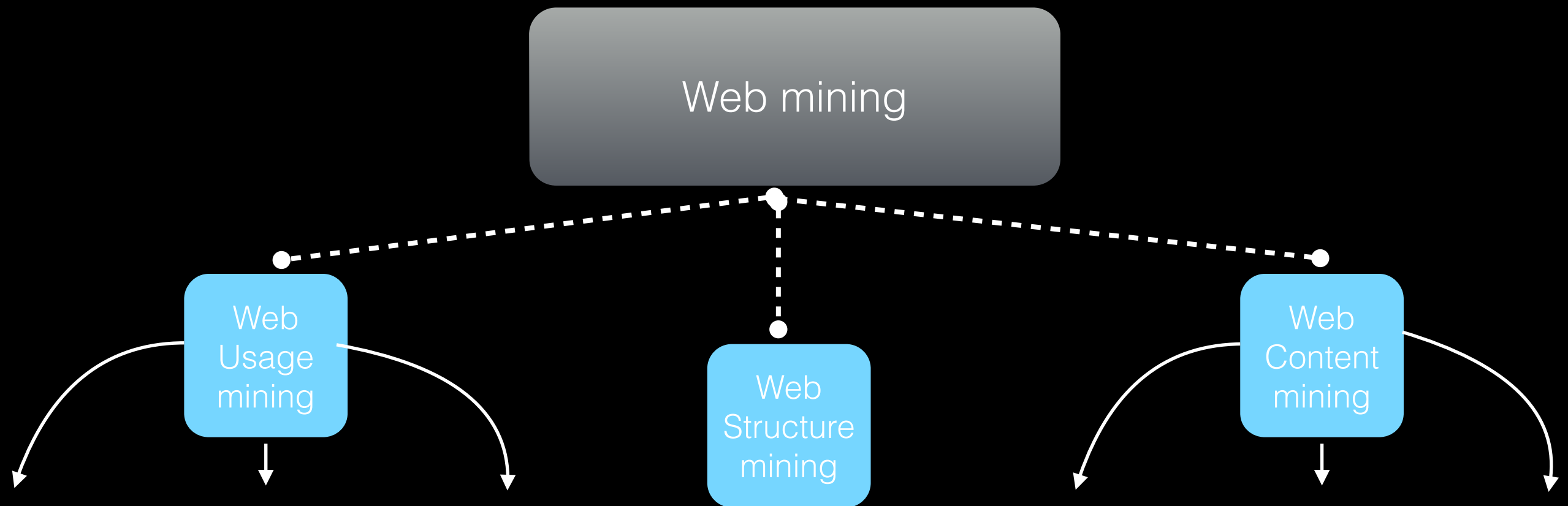
- ▶ Concepts & Définitions
- ▶ Indexation & Crawl
- ▶ Introduction aux techniques de webdatamining
- ▶ Applications a la classification des documents

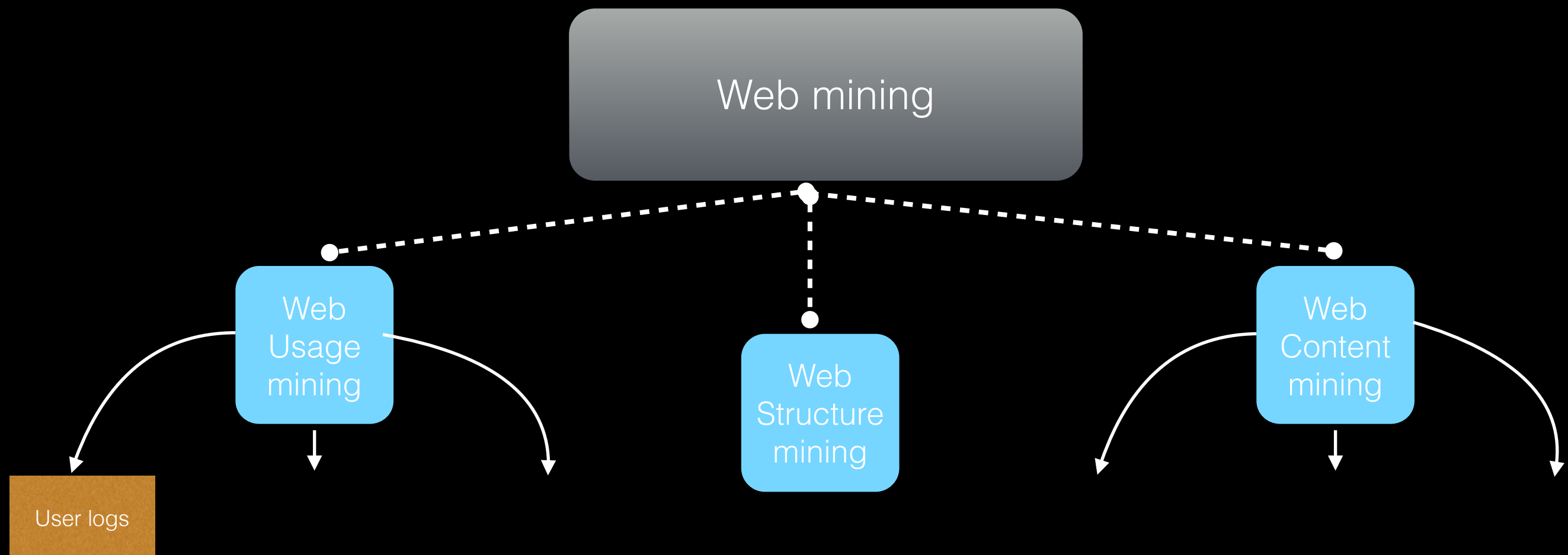


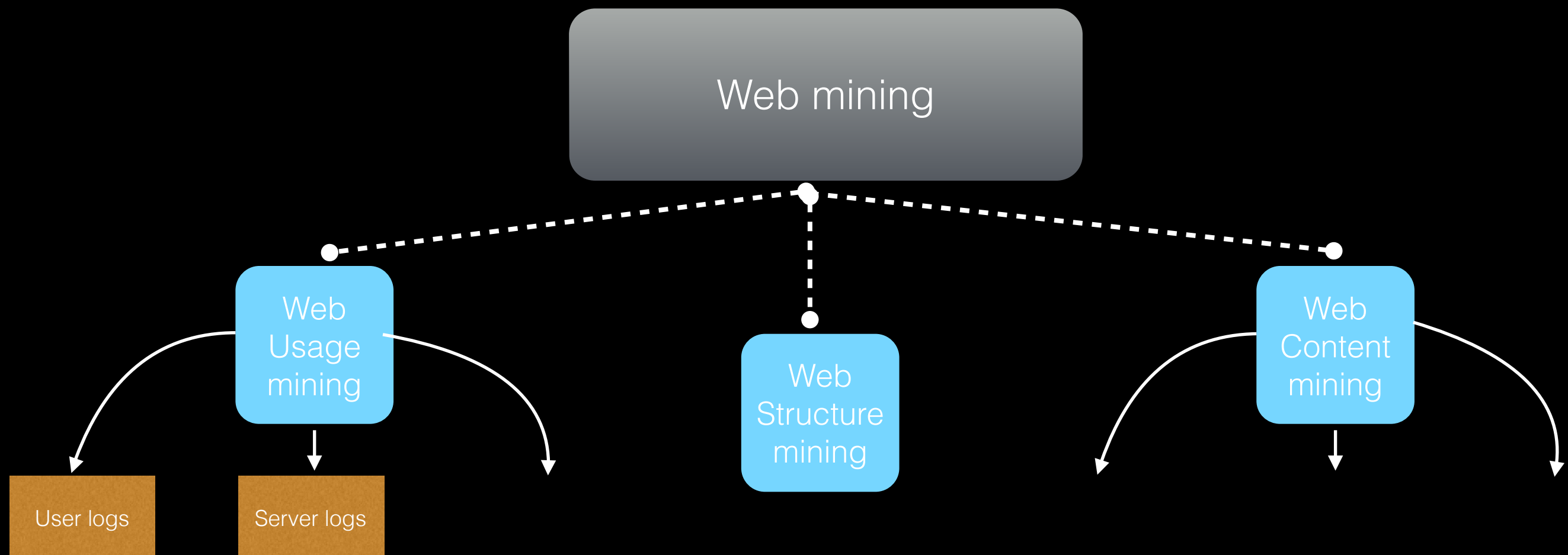


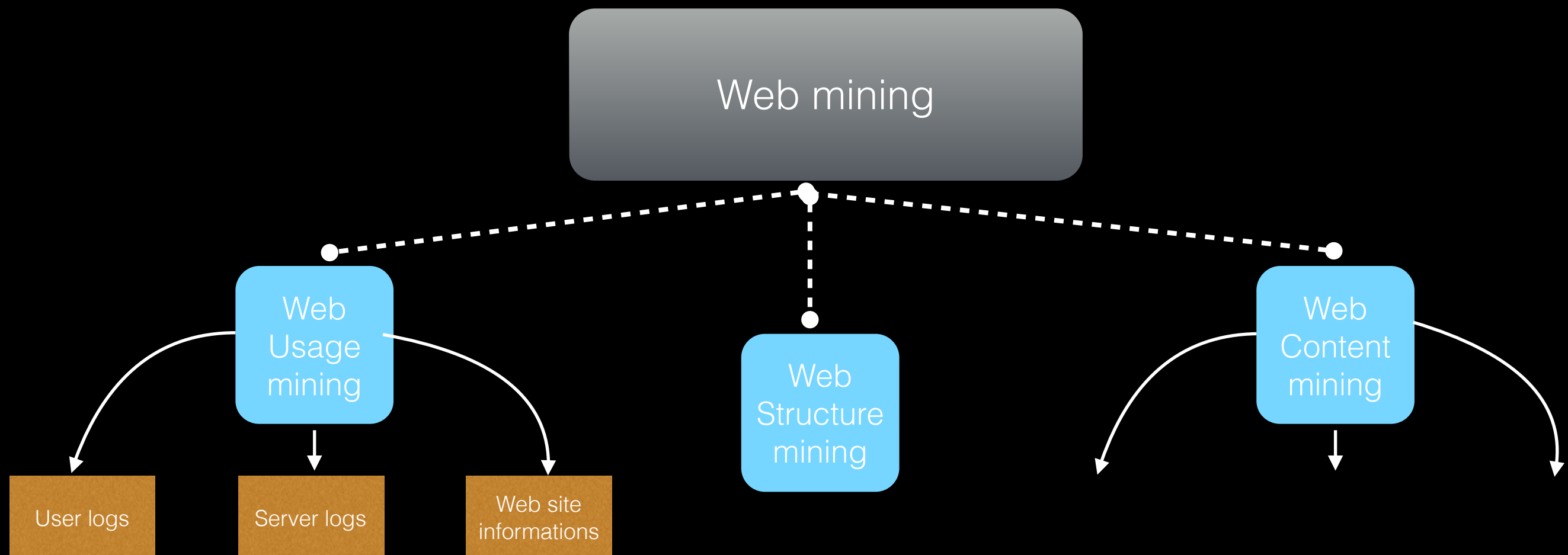


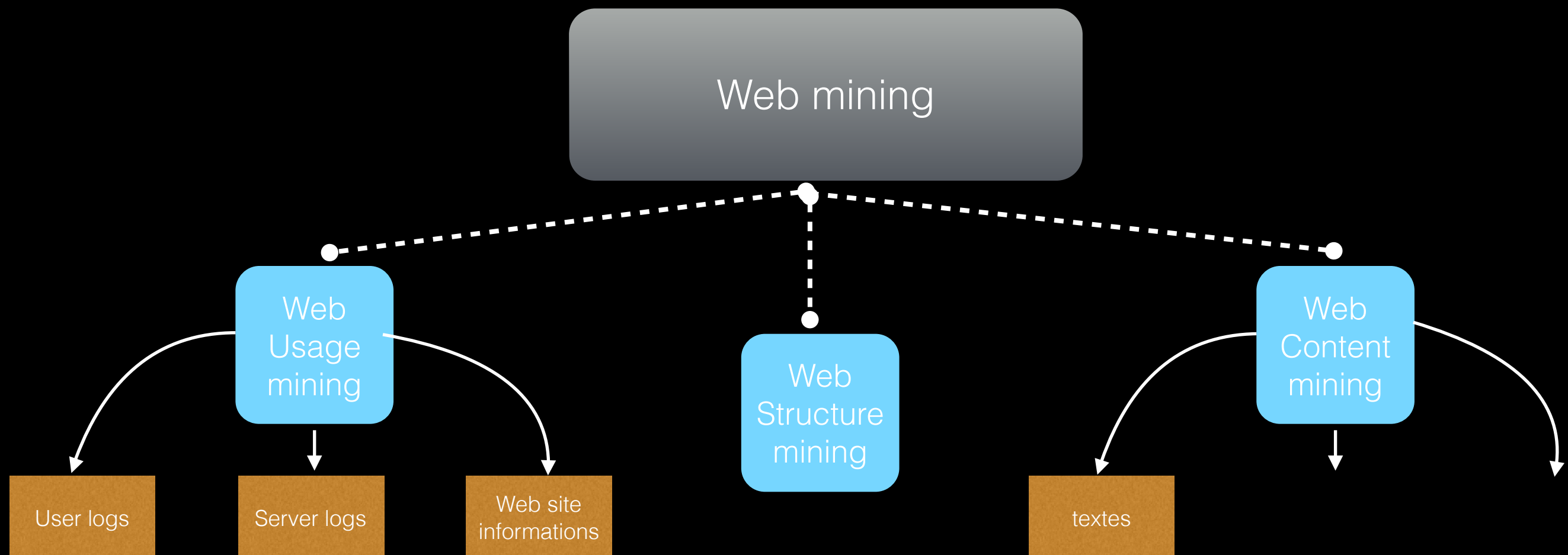


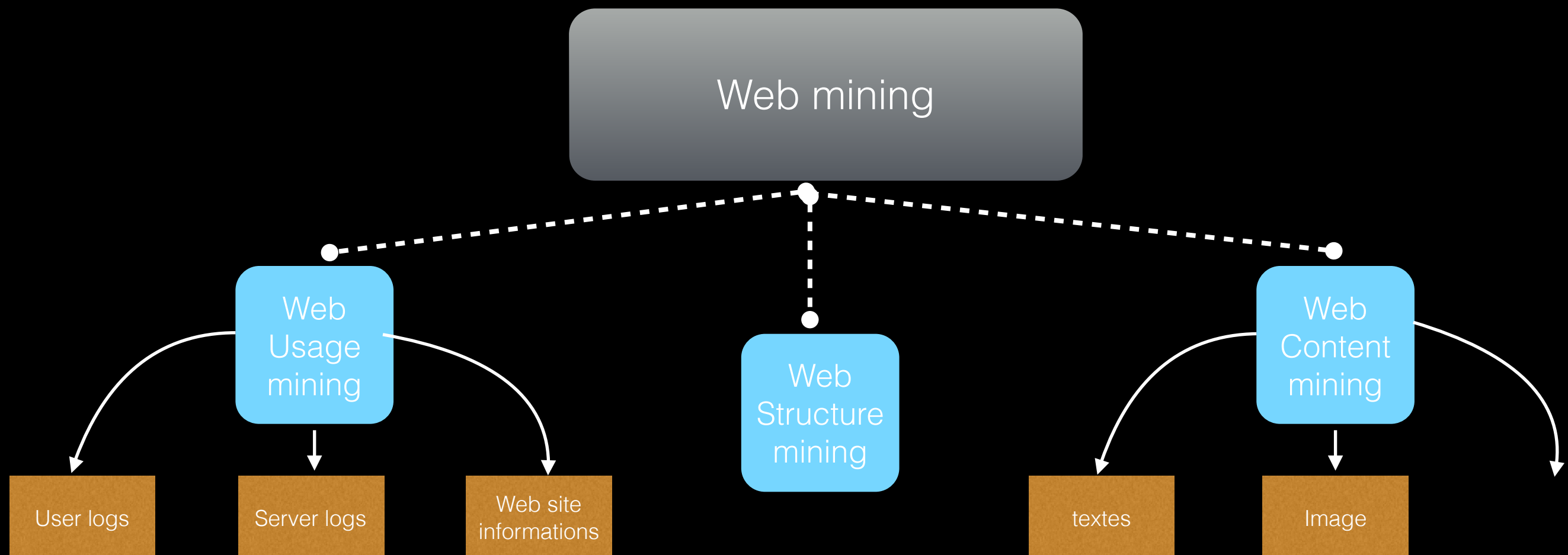


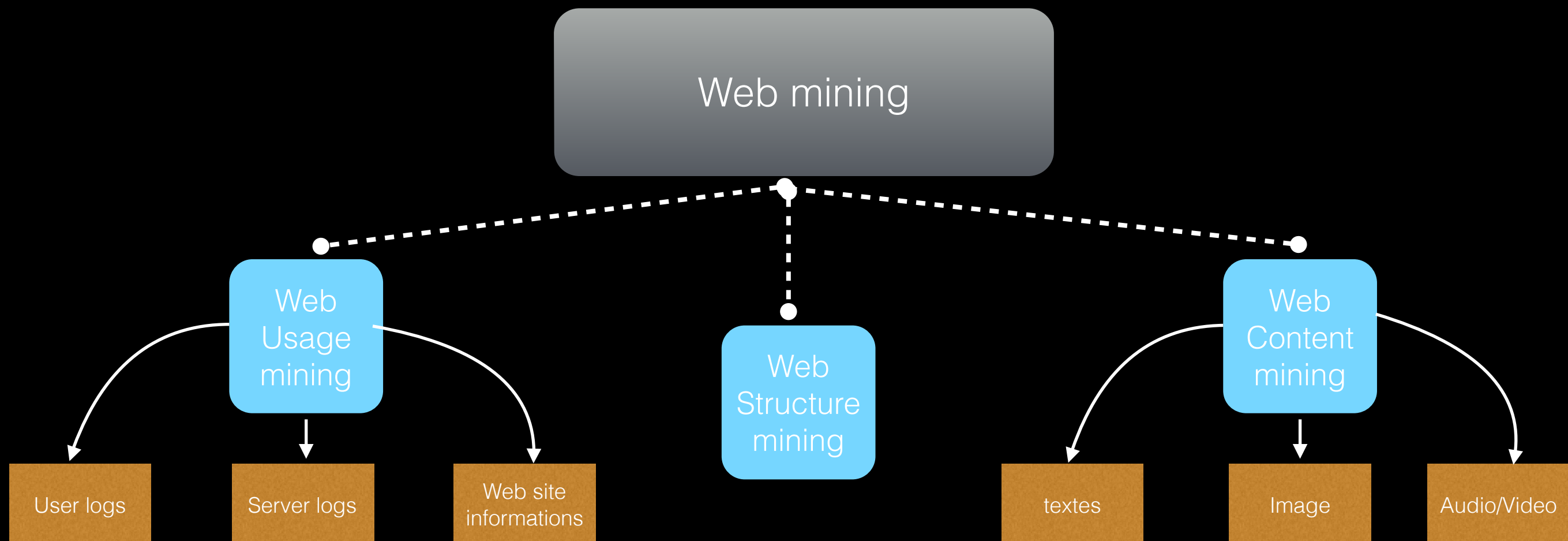


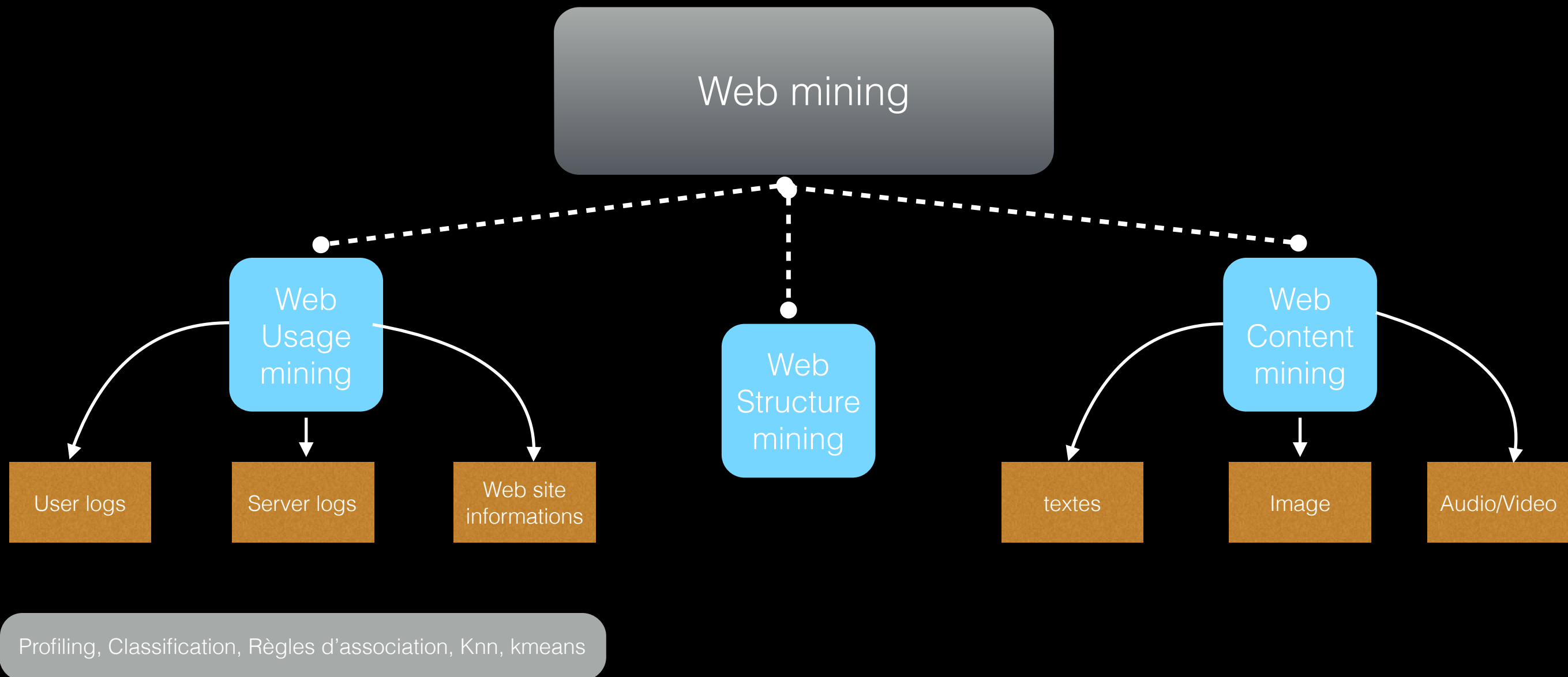


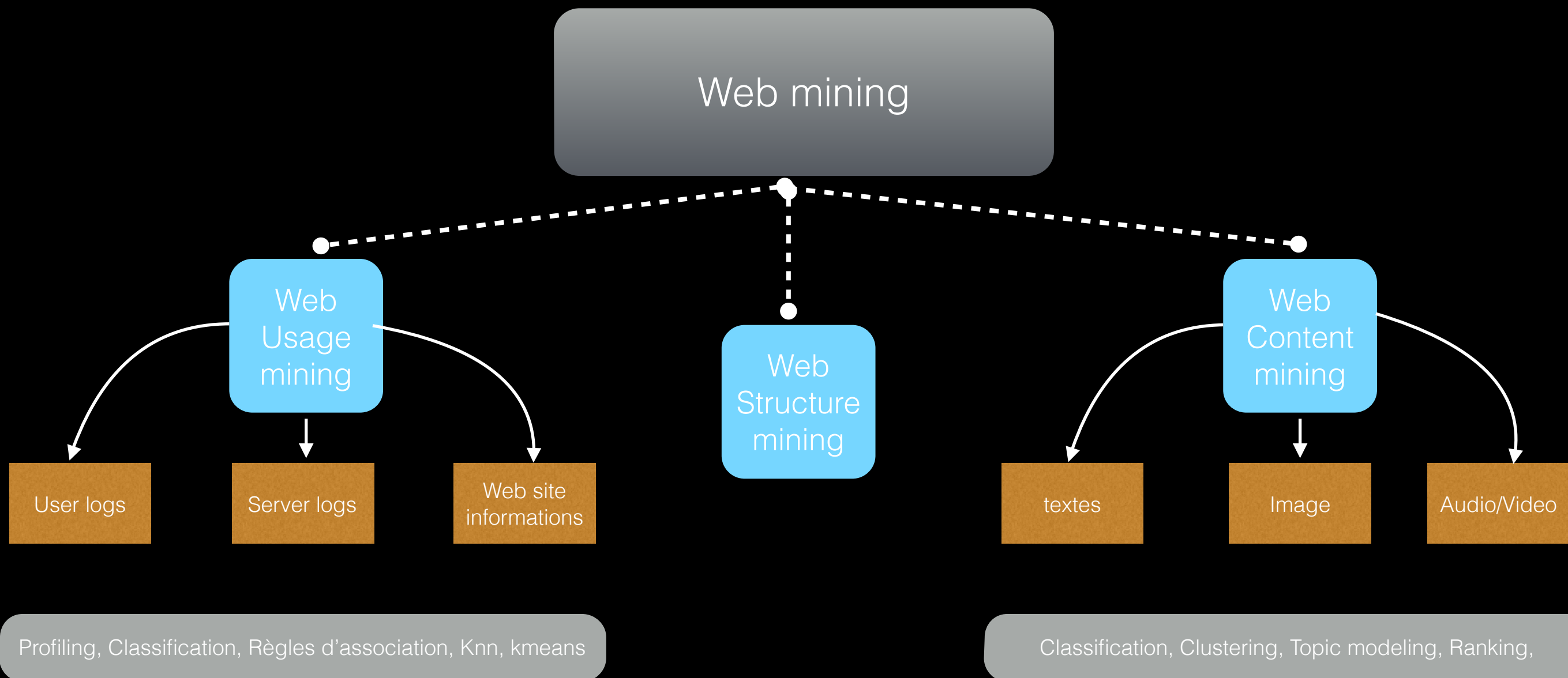




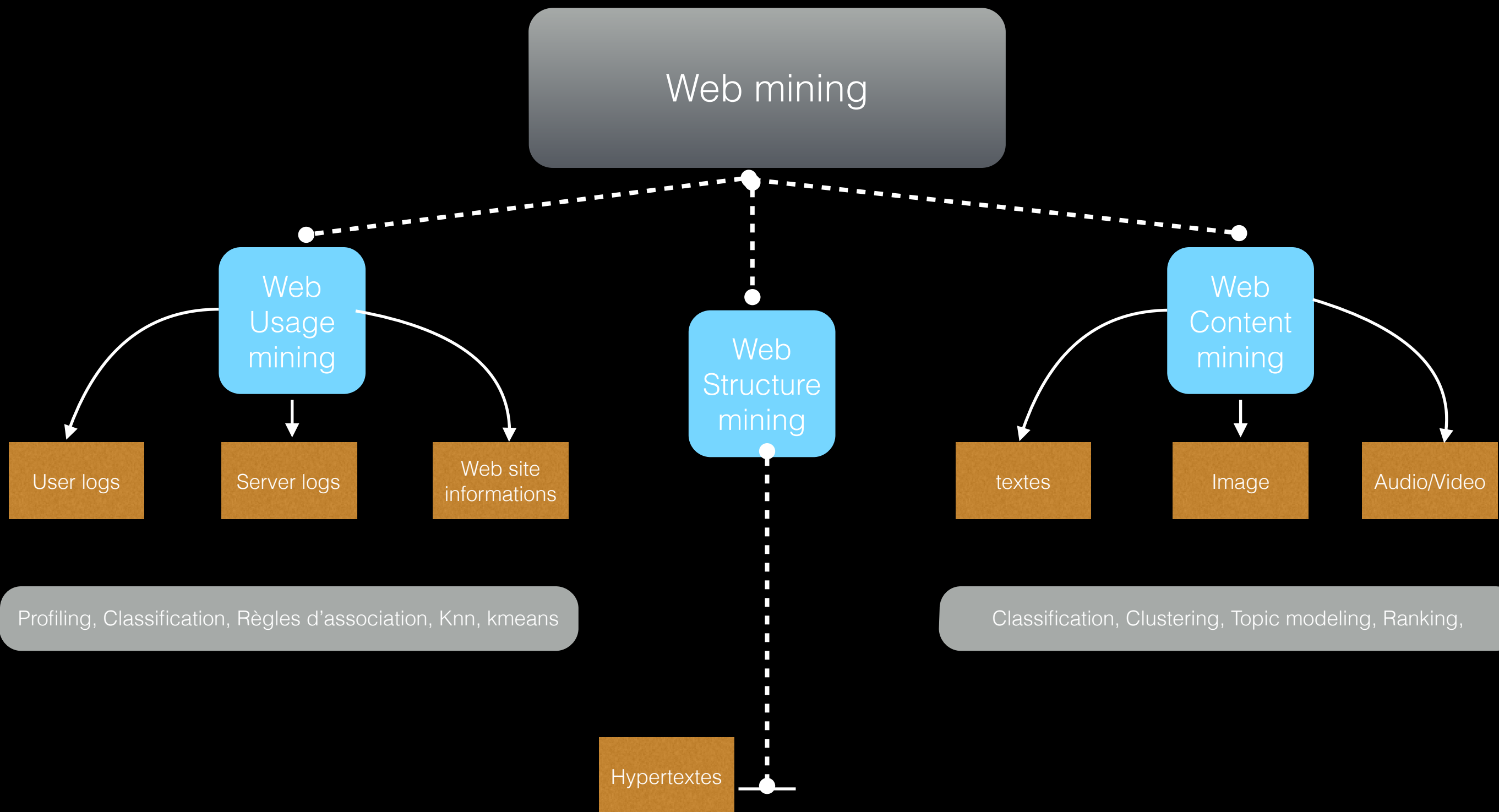


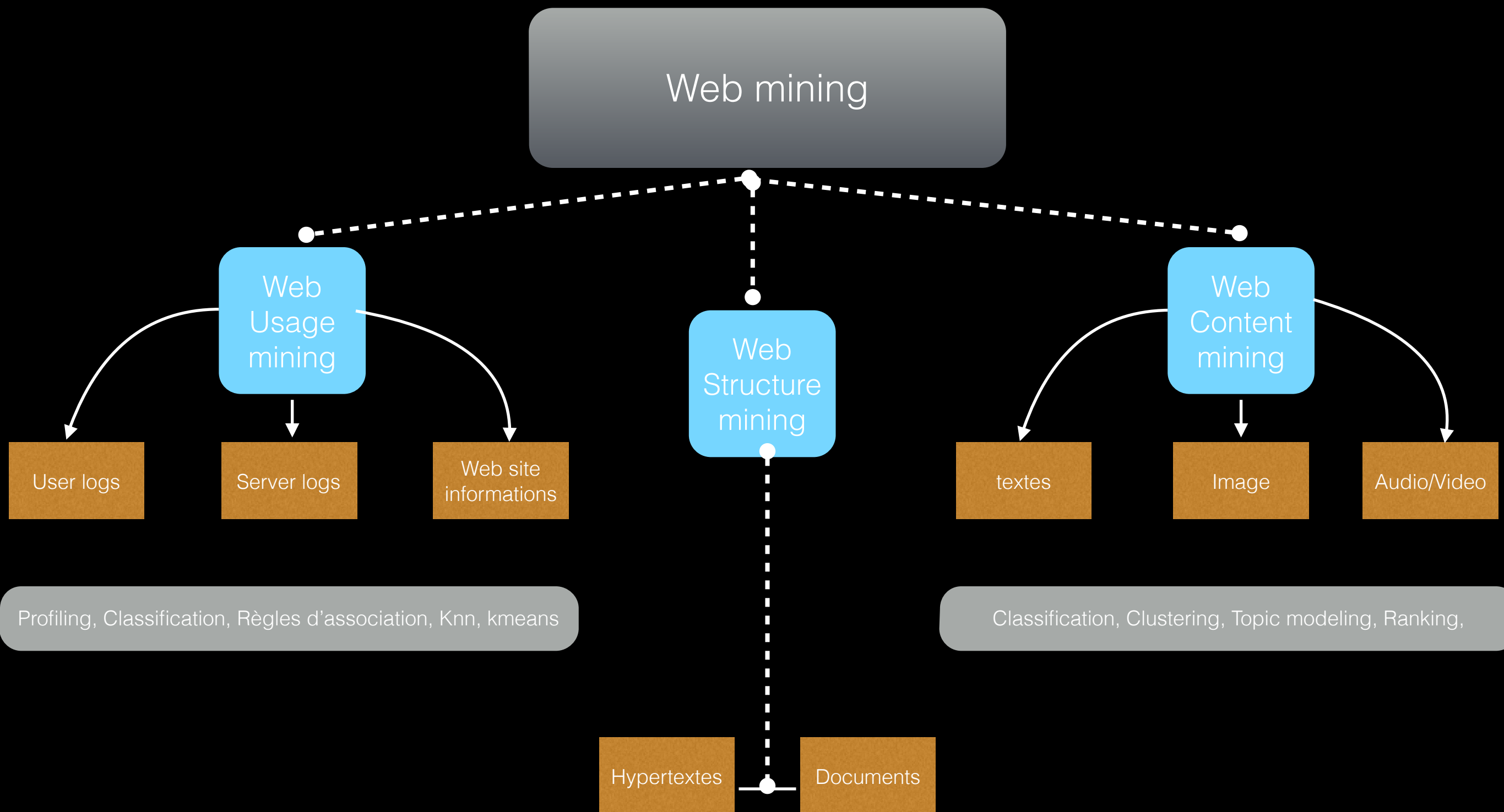


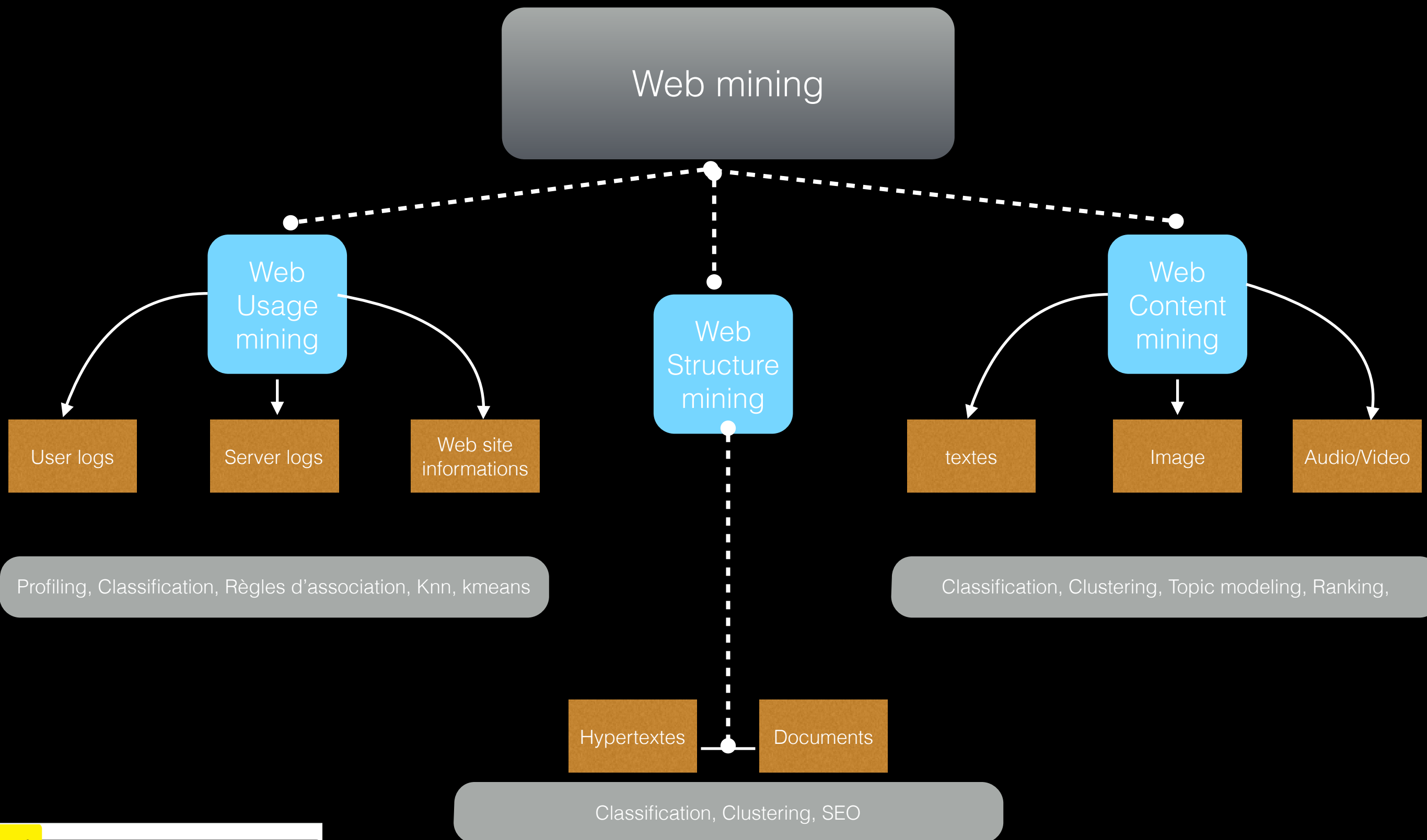




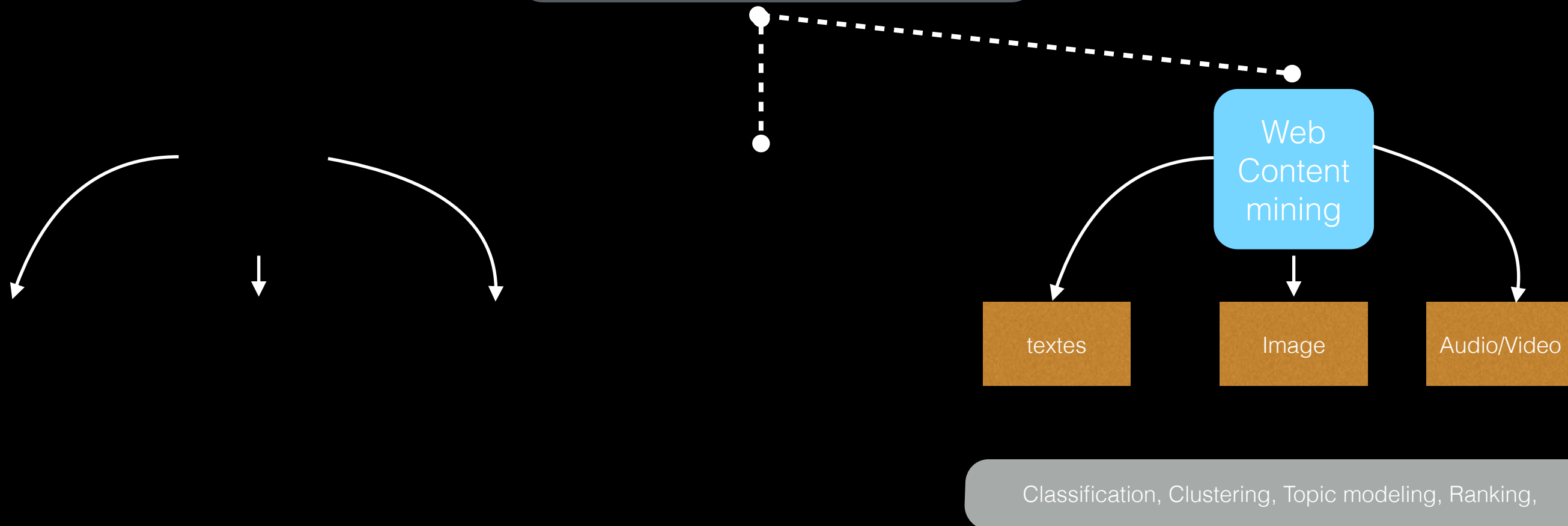




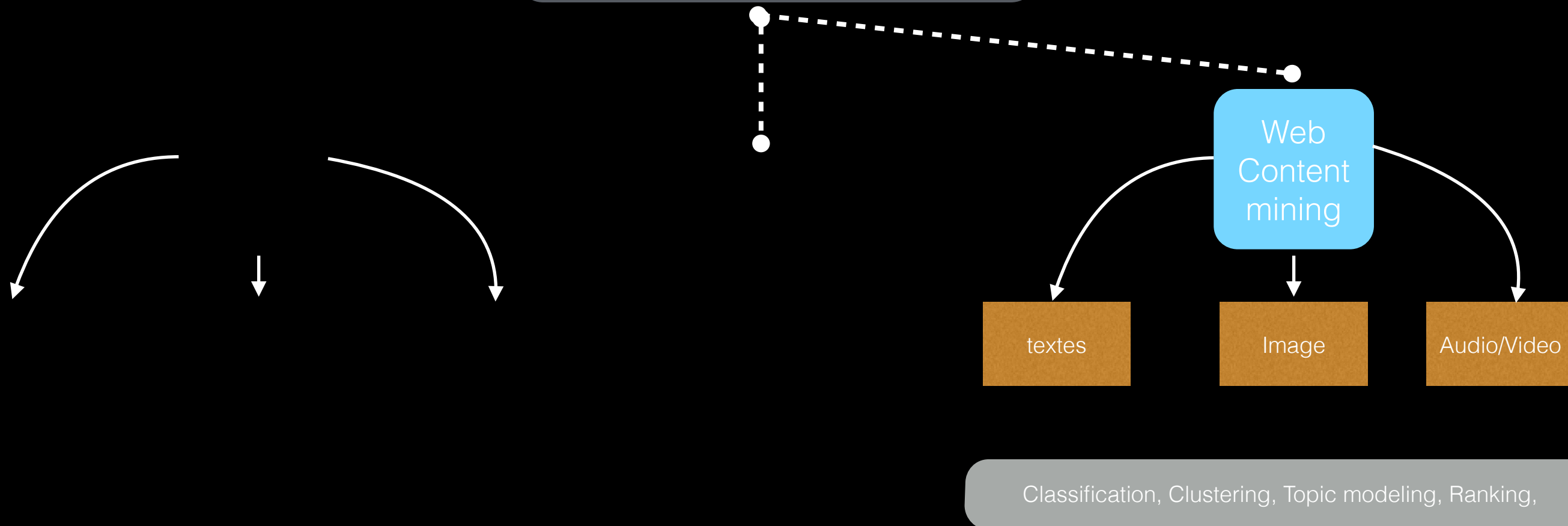


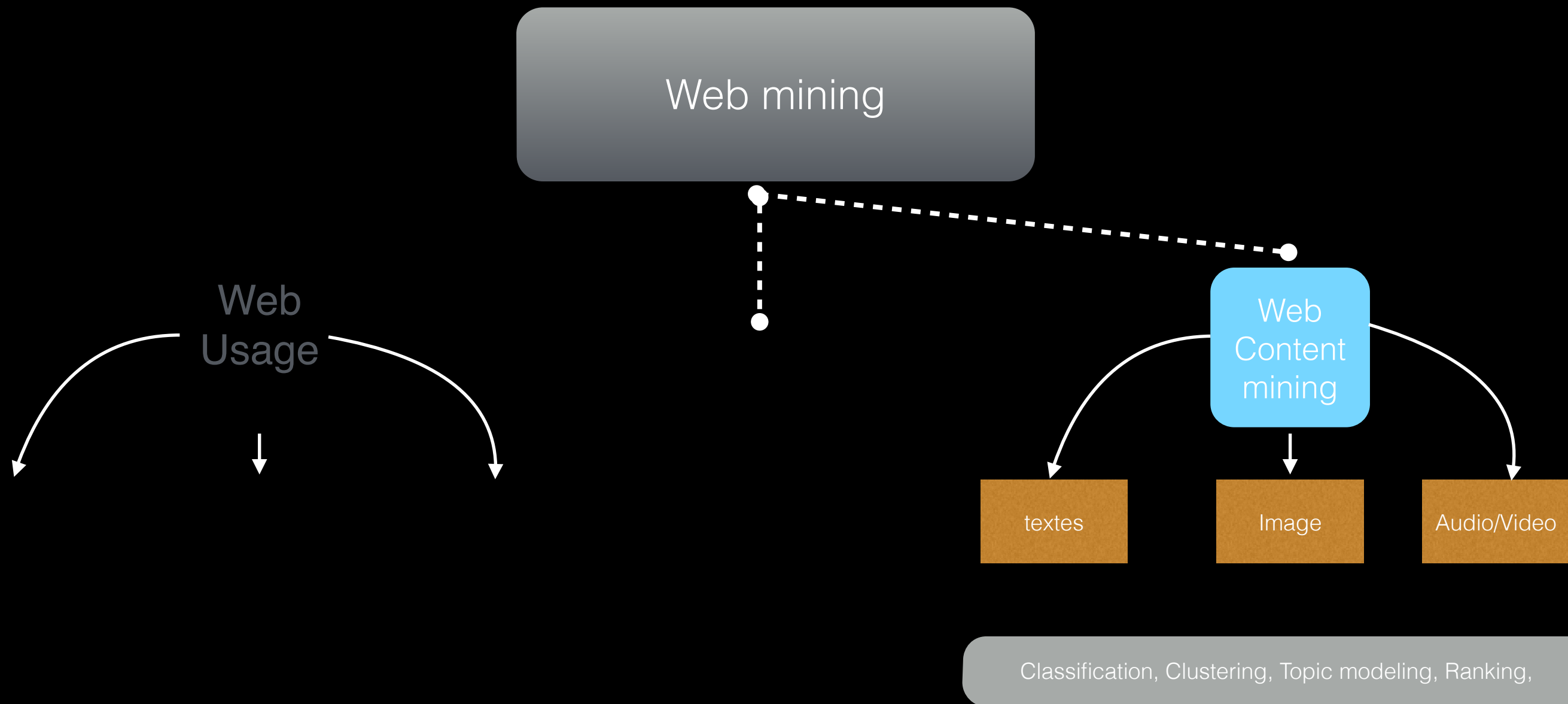


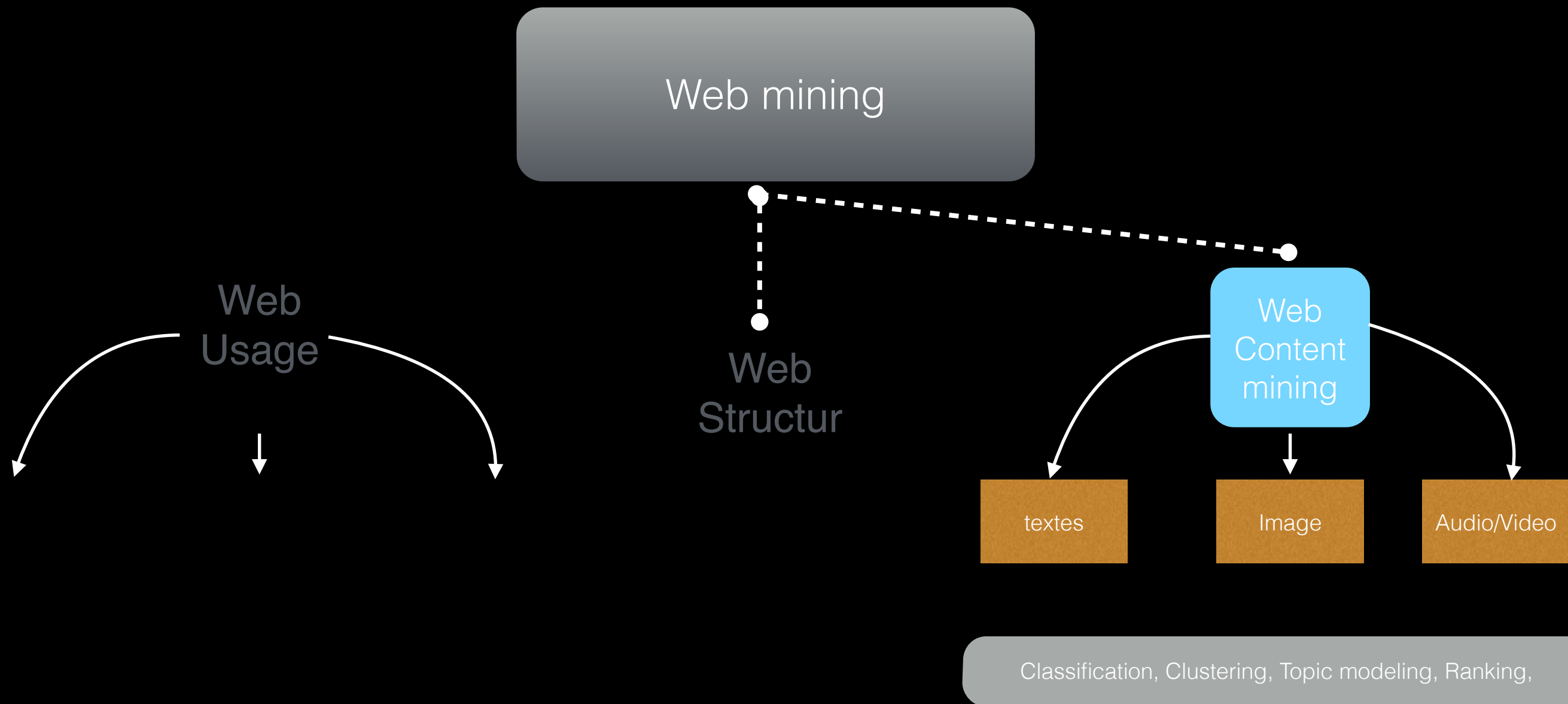
Web mining

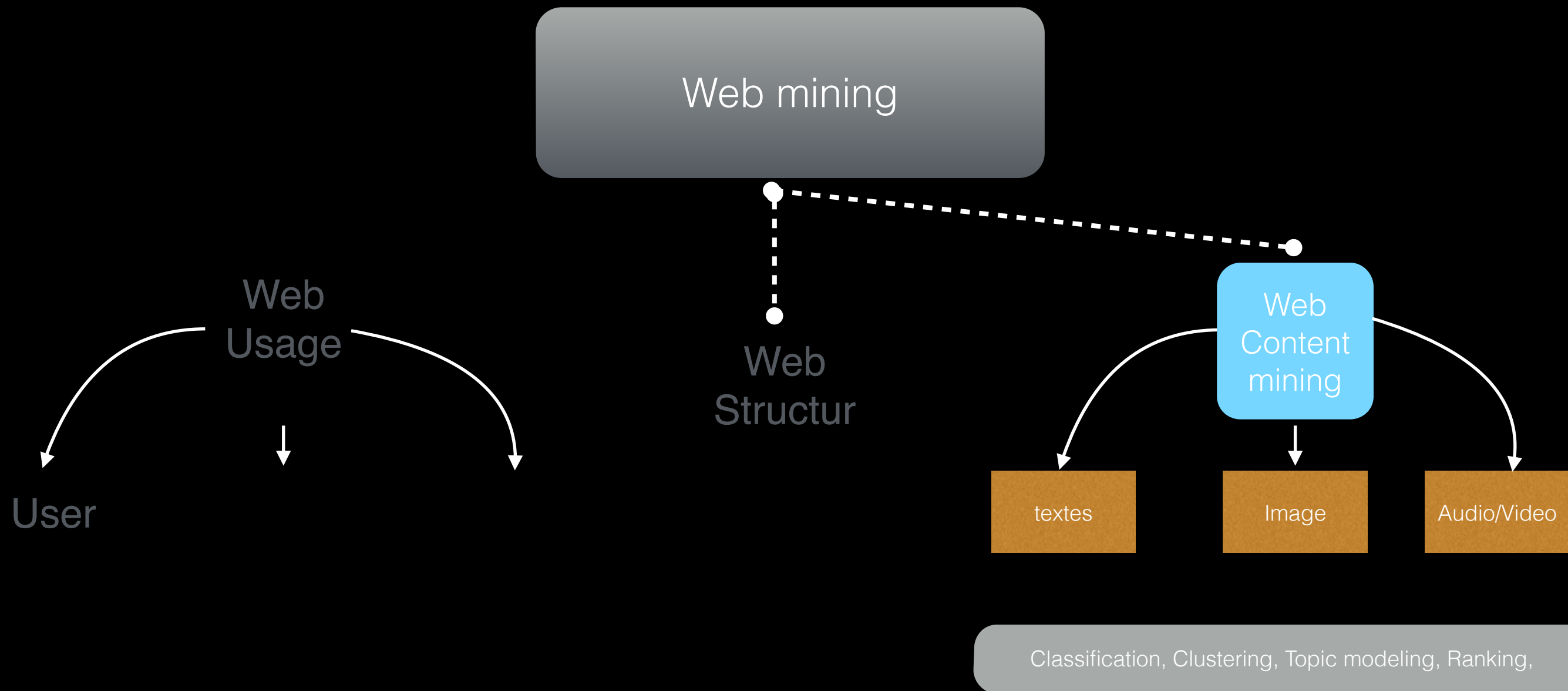


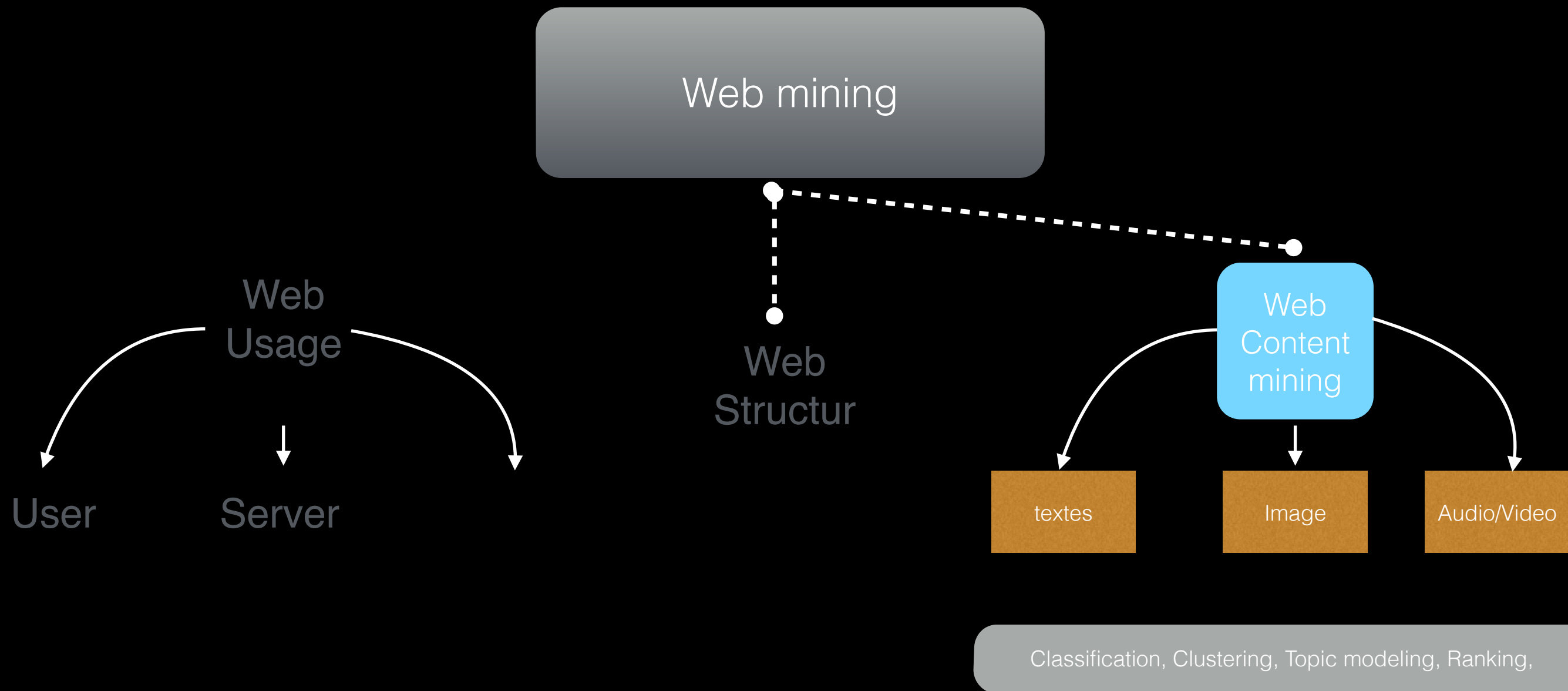
Web mining

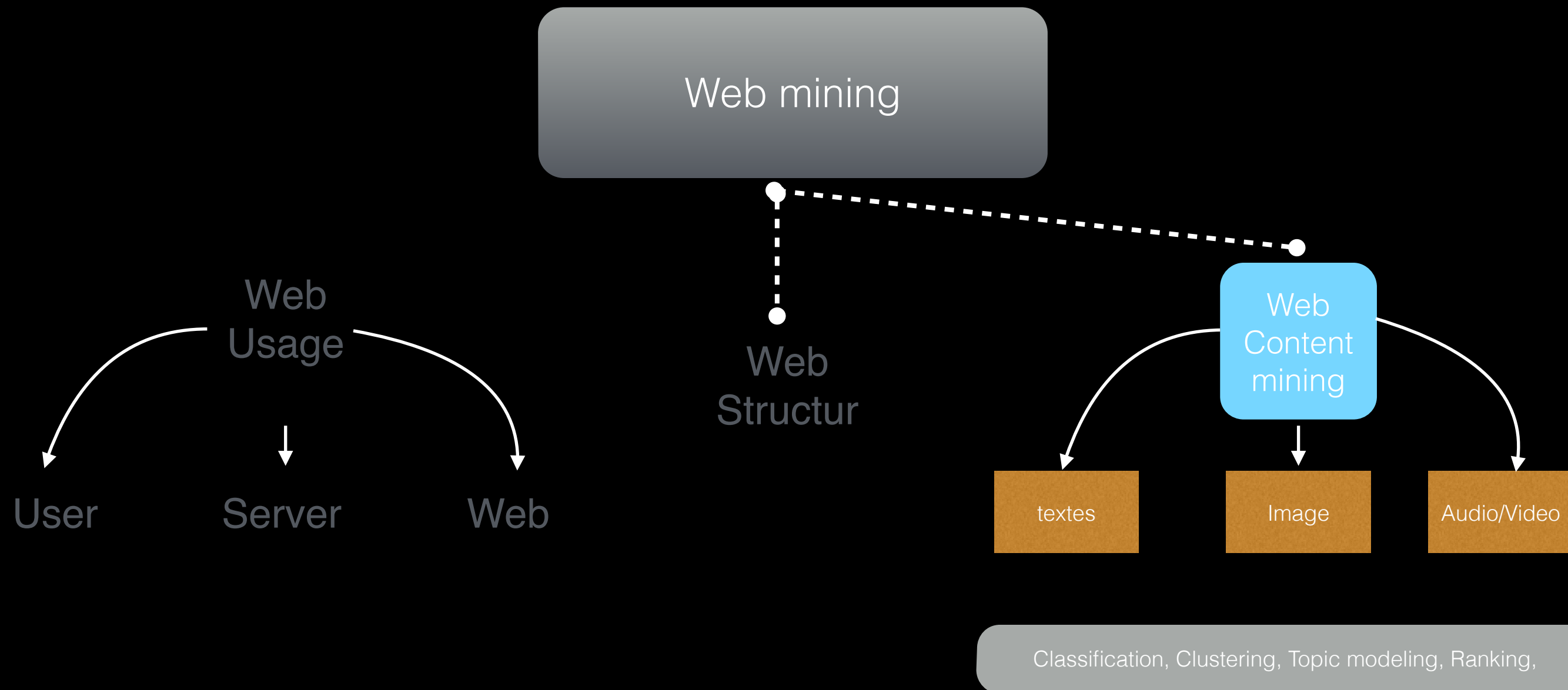


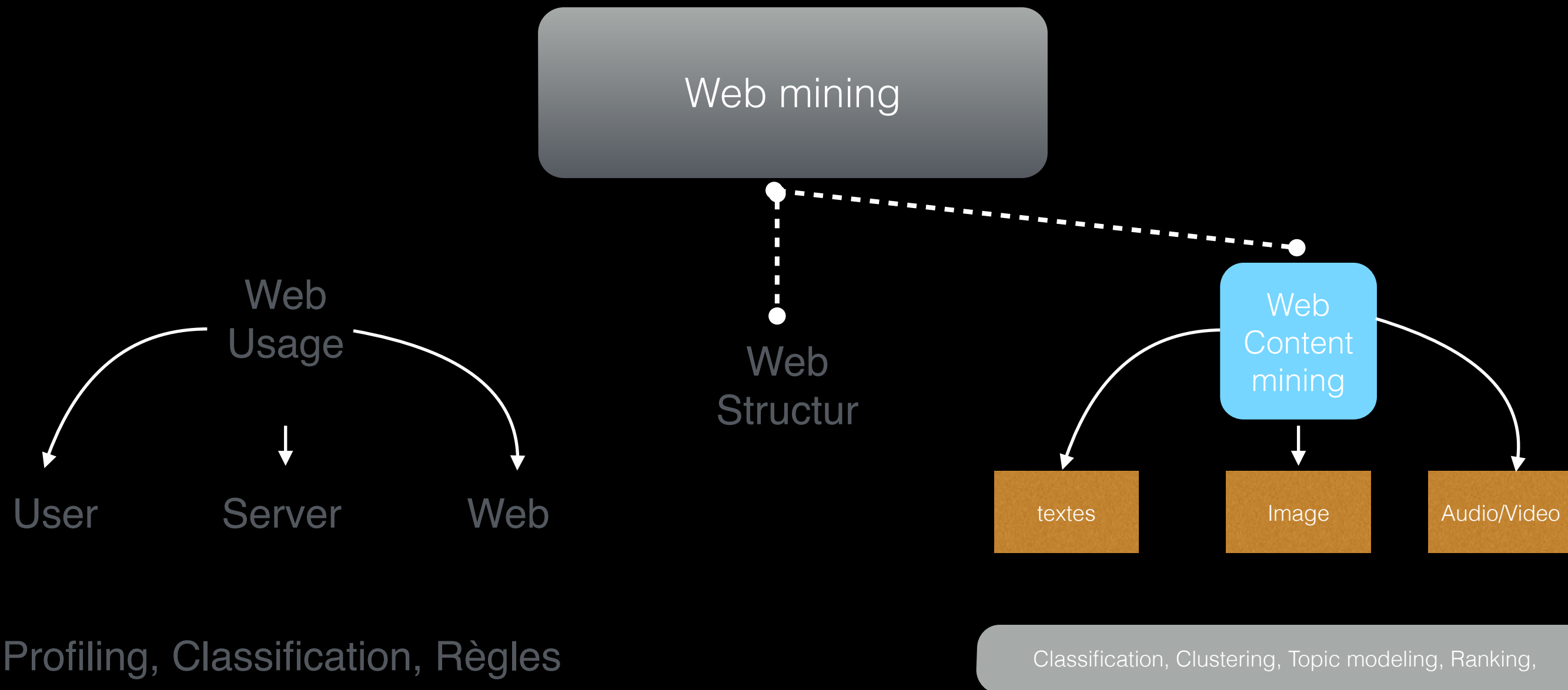


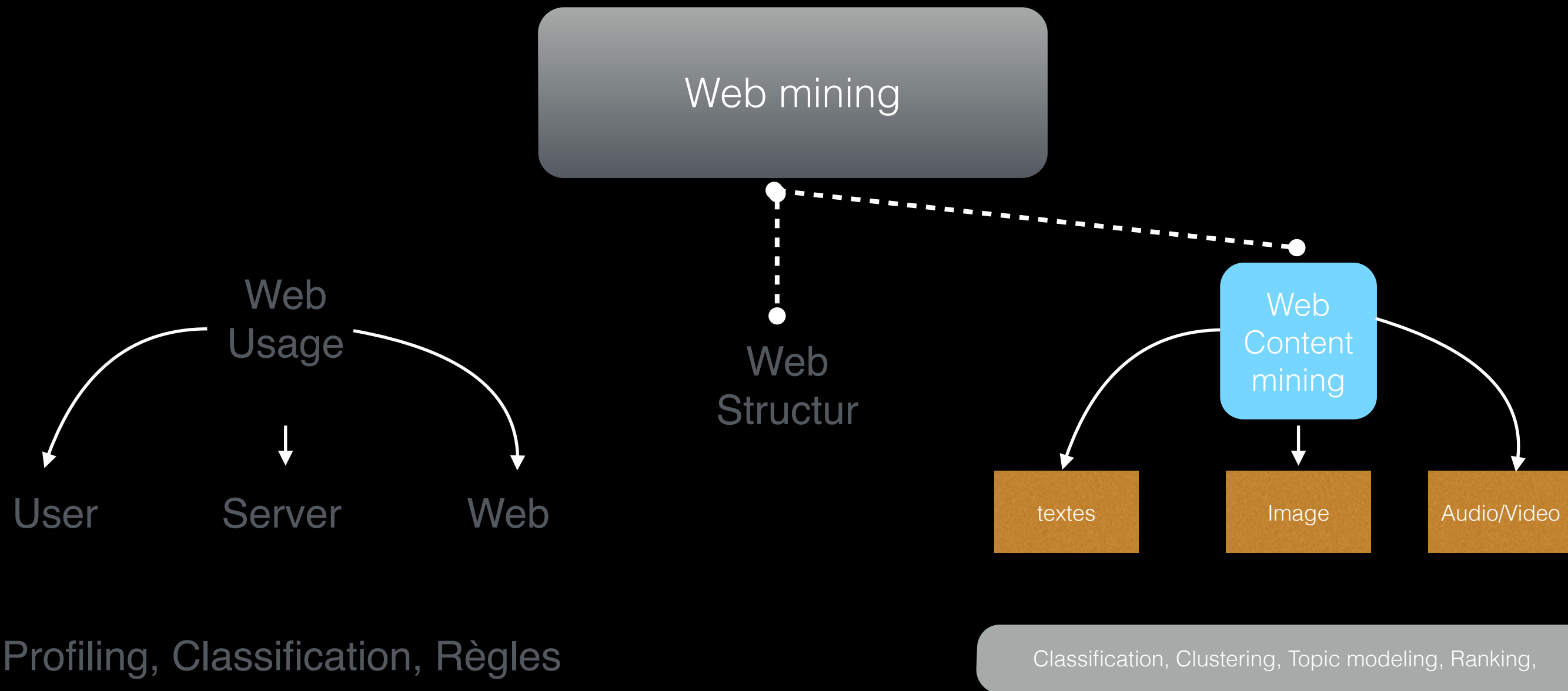


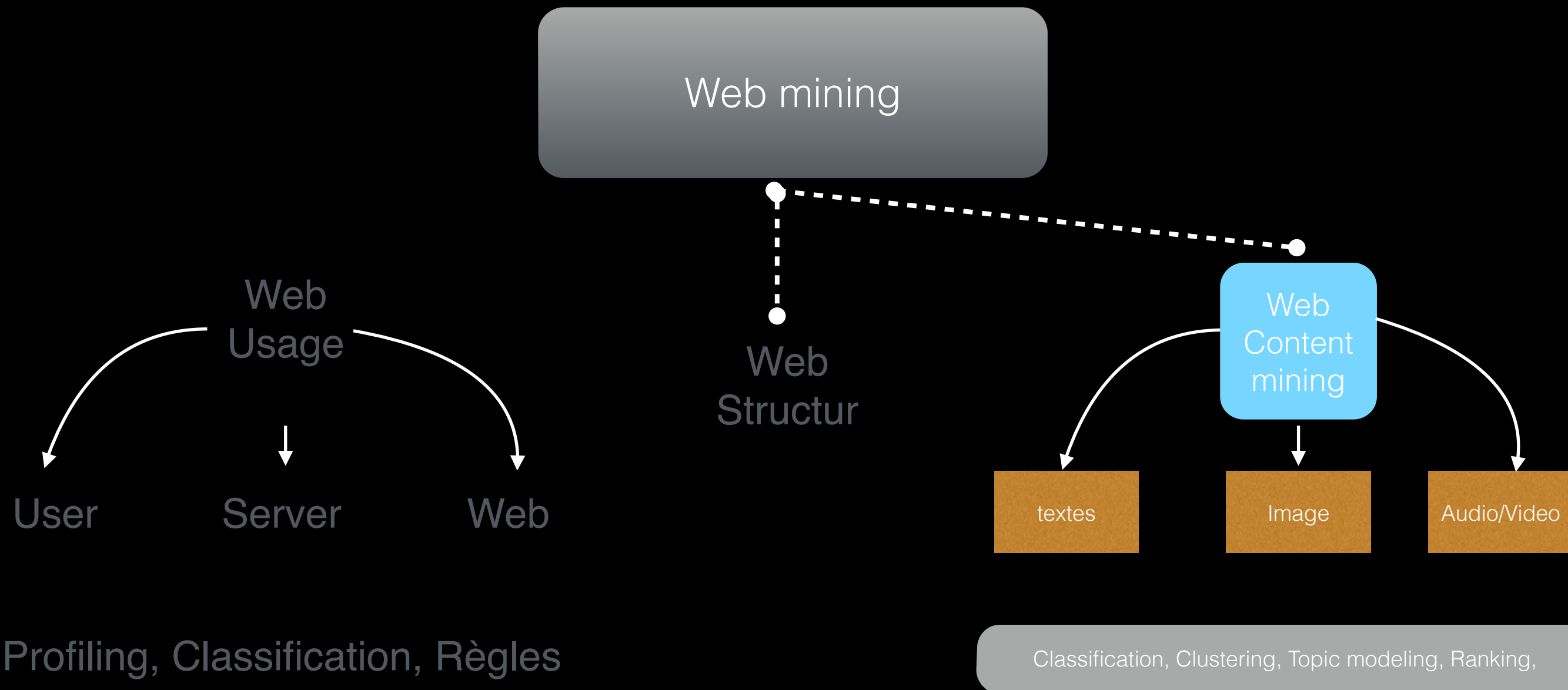




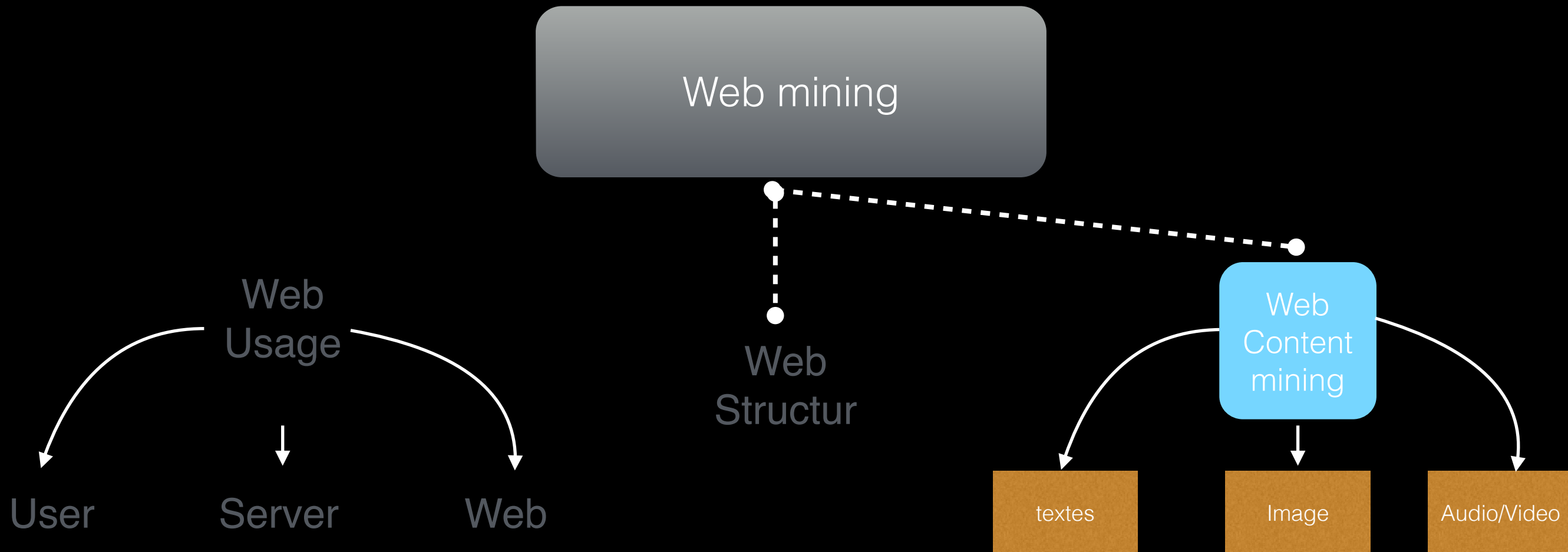








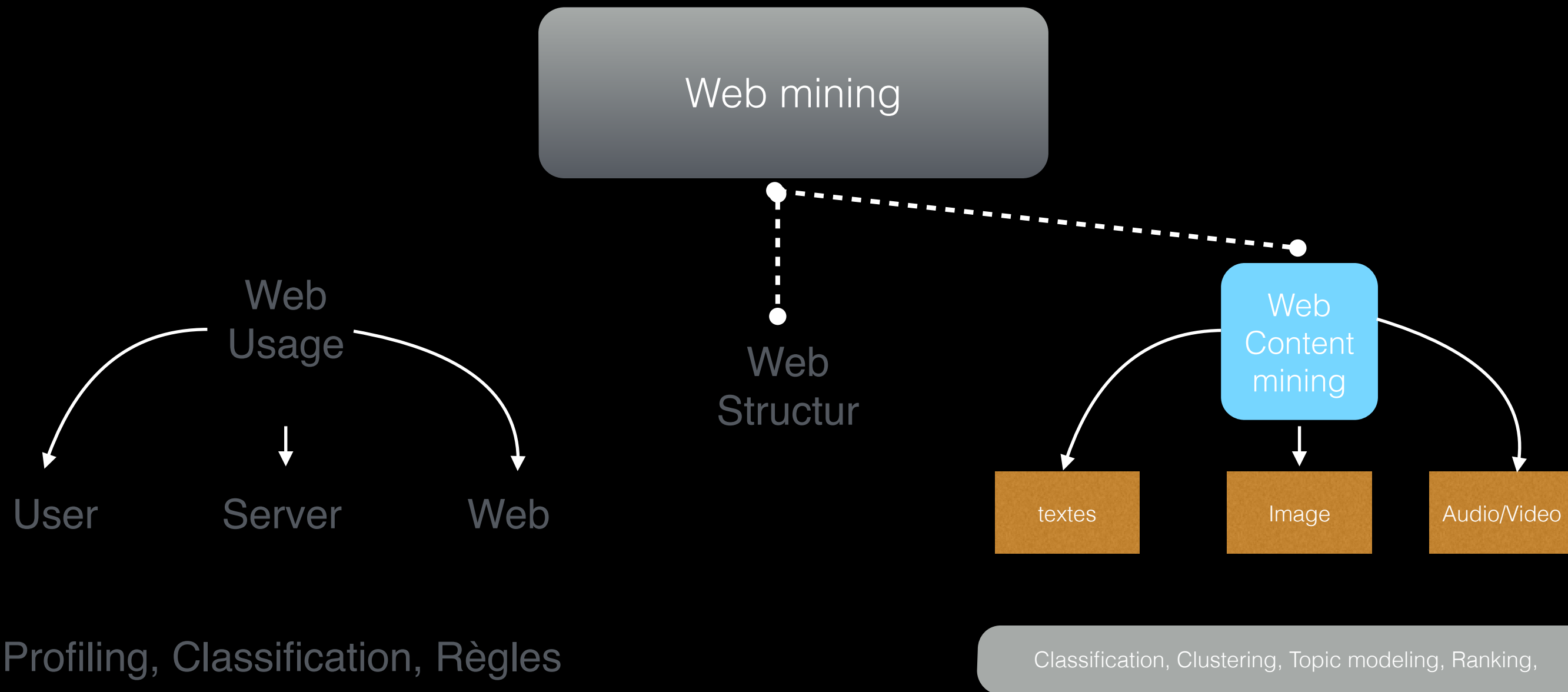
Hypert —



Profiling, Classification, Règles

Classification, Clustering, Topic modeling, Ranking,

Hypert — Docum



Hypert — Docum

Classification, Clustering, SEO

Programme

- ▶ Concepts & Définitions
- ▶ Indexation & Crawl
- ▶ Introduction aux techniques de webdatamining
- ▶ Applications a la classification des documents

Concepts et Définitions

Premières définitions

- La recherche d'informations est le coeur du search

- ▶ La recherche d'informations est le coeur du search
- ▶ Web search has its root in **information retrieval** (or IR for short), a field of study that helps **the user find needed information from a large collection of text documents**. [Liu]

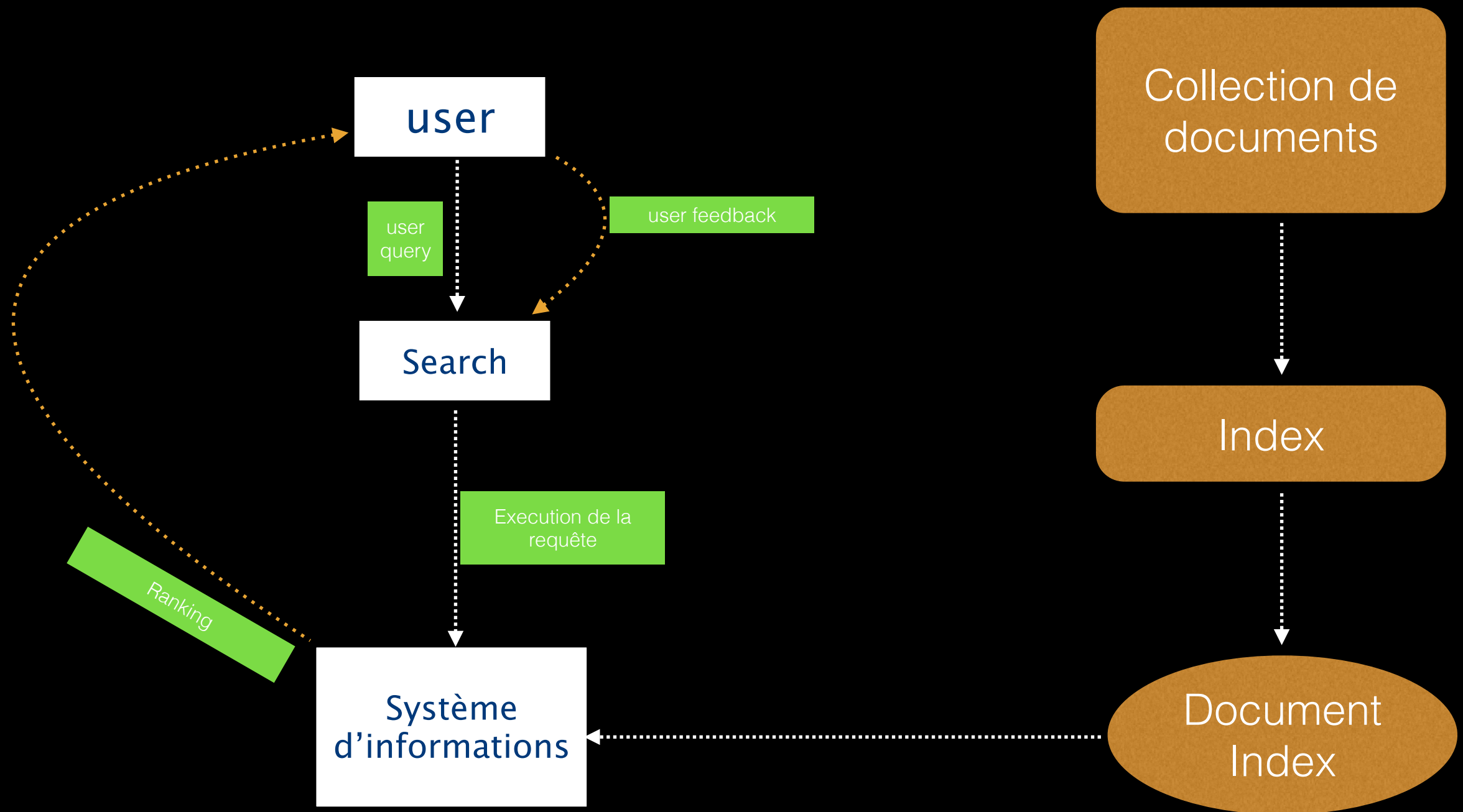
- ▶ La recherche d'informations est le coeur du search
- ▶ Web search has its root in **information retrieval** (or IR for short), a field of study that helps **the user find needed information from a large collection of text documents**. [Liu]
- ▶ en RI => Document

- ▶ La recherche d'informations est le coeur du search
- ▶ Web search has its root in **information retrieval** (or IR for short), a field of study that helps **the user find needed information from a large collection of text documents**. [Liu]
- ▶ en RI => Document
- ▶ Web => Document = page web

- ▶ La recherche d'informations est le coeur du search
- ▶ Web search has its root in **information retrieval** (or IR for short), a field of study that helps **the user find needed information from a large collection of text documents**. [Liu]
- ▶ en RI => Document
- ▶ Web => Document = page web
- ▶ RI = Chercher de l'information pertinente à partir d'une requête

- ▶ La recherche d'informations est le coeur du search
- ▶ Web search has its root in **information retrieval** (or IR for short), a field of study that helps **the user find needed information from a large collection of text documents**. [Liu]
- ▶ en RI => Document
- ▶ Web => Document = page web
- ▶ RI = Chercher de l'information pertinente à partir d'une requête
- ▶ Plusieurs notions en découlent

- ▶ La recherche d'informations est le coeur du search
- ▶ Web search has its root in **information retrieval** (or IR for short), a field of study that helps **the user find needed information from a large collection of text documents**. [Liu]
- ▶ en RI => Document
- ▶ Web => Document = page web
- ▶ RI = Chercher de l'information pertinente à partir d'une requête
- ▶ Plusieurs notions en découlent
 - ➔ Mots clés, Ranking, Pertinence, Base de données, Indexation



- ▶ Recherches par mots clés : Exemple de PagesJaunes
 - ▶ L'utilisateur exprime sa recherche à l'aide de quelques mots clés. Le moteur répond par une liste de documents contenant les mots clés (ex : médecin)
- ▶ Recherche full text ou de documents : Exemple de google
 - ▶ L'utilisateur exprime sa recherche par une phrase et le moteur remonte des documents qui contiennent au moins une fois la phrase
- ▶ Recherche approchée :
 - ▶ Comme la recherche full, mais sont remontés les documents qui contiennent une partie seulement de la phrase sont remontées
- ▶ Recherche en langage naturel
 - ▶ Les opérations de recherche peuvent être très complexes

Programme

- ▶ Concepts & Définitions
- ▶ Indexation & Crawl
- ▶ Introduction aux techniques de webdatamining
- ▶ Applications a la classification des documents

- ▶ Un index est un **module qui enregistre** les documents bruts et permet ainsi une recherche efficace
- ▶ L'action de rechercher l'information consiste donc à :
 - ▶ Calculer le degré de pertinence de chaque document **indexé** avec requête de l'utilisateur
 - ▶ Ordonnancer une liste de documents pertinents à renvoyer l'utilisateur

► Modèle Booléen

In the Boolean model, documents and queries are represented as sets of terms. That is, each term is only considered present or absent in a document. [LIU]

La recherche dans un modèle booléen s'effectue par l'utilisation de l'algèbre des prédicats (et, ou, non)

Le système d'information renvoie des documents pour lesquels la requête est TRUE

Mot clé à retenir : **EXACT MATCH**

► Modèle à vecteurs espaces-états

A document in the vector space model is represented as a weight vector, in which each component weight is computed based on some variation of TF or TF-IDF scheme.[LIU]

TF : Terms Frequency (Nombre de fois où le terme apparaît dans le document)

IDF : Inverse Document Frequency (Nombre de documents dans lequel le terme apparaît)

La recherche dans le système d'informations se fait sur la base de mesure de similarité entre documents => Cosinus, Indice de Jaccard, etc...

Programme

- ▶ Concepts & Définitions
- ▶ Indexation & Crawl
- ▶ Introduction aux techniques de webdatamining
- ▶ Applications a la classification des documents

► TDM ou DTM

Une matrice de matrice document–terme ou terme–document est une matrice au sens mathématique du terme qui décrit la fréquence des termes qui apparaissent dans une collection de documents.

Dans une tdm(term–document–matrix), les lignes correspondent aux documents de la collection et les colonnes correspondent à des termes. Il existe différents systèmes pour déterminer la valeur que chaque entrée de la matrice doit prendre. le plus répandu est la mesure tf–idf(Term Frequency–Inverse Document Frequency), mais aussi, tf(Term Frequency).

PAUSE



Web Datamining

Information Retrieval and Machine Learning for the Web

Guibert J. TCHINDE

ENSAI | Solocal group – 2014