

# Algèbres de collections

E.Coquery

`emmanuel.coquery@univ-lyon1.fr`

`https:`

`//forge.univ-lyon1.fr/mif24-bdnosql/mif24-bdnosql`

# Collections

"Paquets" de valeurs :

- Avec ou sans doublons
- Triés ou non
- Contenant des valeurs homogènes

# Des opérations sur les collections communes à de nombreux environnements

Lesquelles ?

- Transformer
- Filtrer, extraire
- Combiner
- Agréger
- Trier

Où ?

- Dans les bases de données (SQL, NoSQL)
- Dans les langages de programmation (Javascript, Python, Java, OCaml, etc)

# Algèbre

Dans le cadre du cours :

- Ensemble : collections ordonnées de documents JSON / valeurs
  - Le nombre d'occurrences compte
  - L'ordre d'apparition des éléments compte
- Opérations :
  - $Map_f$
  - $Filter_f$
  - $Join_f$
  - $Agg_{f_m, f_a, k}, FlatMap_f$
  - $Sort_f$

# Mini-langage pour exprimer les calculs

$$\begin{aligned}
 d &:= & 1 & \mid & 2 & \mid & \dots \\
 & \mid & 1.0 & \mid & 5.7 & \mid & \dots \\
 & \mid & \text{"truc"} & \mid & \dots \\
 & \mid & d(d) & \mid & d_1 \text{ op } d_2 \\
 & \mid & \lambda x. d \\
 & \mid & \{a_1 : d_1, \dots, a_n : d_n\} \\
 & \mid & d.a \\
 & \mid & [] & \mid & [d] \\
 \text{op} &:= & + & \mid & * & \mid & = & \mid & ++ & \mid & \dots
 \end{aligned}$$

Rmq :  $d_1 \text{ op } d_2$  est une écriture pour  $\text{op}(d_1)(d_2)$

# Types des constantes

$$(\text{Const}) \frac{}{\Gamma \vdash c : \tau}$$

en prenant  $\tau$  et  $c$  comme suit :

- type *int* : 1 , 2 , ...
- type *float* : 1.0 , 0.3 , 10.42 , ...
- type *string* : "truc" , ...

# Typage des fonctions

$$(\text{App}) \frac{\Gamma \vdash f : \tau \rightarrow \tau' \quad \Gamma \vdash d : \tau}{\Gamma \vdash f(d) : \tau'}$$

$$(\text{Lambda}) \frac{\Gamma ++ [x : \tau] \vdash d : \tau'}{\Gamma \vdash \lambda x. d : \tau \rightarrow \tau'}$$

$$(\text{Var}) \frac{}{\Gamma \vdash x : \tau} \text{ si } x : \tau \in \Gamma$$

# Typage des records, des listes

$$\text{(Field)} \frac{\Gamma \vdash d : \langle a : \tau \rangle}{\Gamma \vdash d.a : \tau}$$

$$\text{(Record)} \frac{\Gamma \vdash d_1 : \tau_1 \quad \dots \quad \Gamma \vdash d_n : \tau_n}{\Gamma \vdash \{a_1 : d_1, \dots, a_n : d_n\} : \langle a_1 : \tau_1, \dots, a_n : \tau_n \rangle}$$

$$\text{(Singleton)} \frac{\Gamma \vdash d : \tau}{\Gamma \vdash [d] : [\tau]}$$

$$\text{(Empty)} \frac{}{\Gamma \vdash [] : [\tau]}$$



# Sous-typage

$$(\text{Sous-typage}) \frac{\Gamma \vdash d : \tau \quad \tau \preceq \tau'}{\Gamma \vdash d : \tau'}$$

## Sous-types (Rappel)

$$(\text{Refl}) \frac{}{\tau \preceq \tau}$$

$$(\text{Trans}) \frac{\tau \preceq \tau' \quad \tau' \preceq \tau''}{\tau \preceq \tau''}$$

$$(\text{AddField}) \frac{}{\langle a_1 : \tau_1, \dots, a_k : \tau_k, a_{k+1} : \tau_{k+1} \rangle \preceq \langle a_1 : \tau_1, \dots, a_k : \tau_k \rangle}$$

$$(\text{STField}) \frac{\tau_k \preceq \tau'_k}{\langle a_1 : \tau_1, \dots, a_k : \tau_k \rangle \preceq \langle a_1 : \tau_1, \dots, a_k : \tau'_k \rangle}$$

$$(\text{STArray}) \frac{\tau \preceq \tau'}{[\tau] \preceq [\tau']}$$

$$(\text{STDict}) \frac{\tau \preceq \tau'}{\{\tau\} \preceq \{\tau'\}}$$

# Fonctions pures

Pure : Résultat de la fonction ne dépend que de ses arguments, pas du contexte.

Dans la suite du cours, on supposera que toutes les fonctions qui paramètrent des opérateurs sont pures.

# Collections

Contenu :

- records, qui ne contiennent que des données
  - pas de fonction
  - pas de variable
- tous les éléments ont le même type  
(mais sous-typage autorisé)

# Algèbre

Opérateurs paramétrés par :

- des fonctions pures
- des noms de champ

$Map_f$ ,  $Filter_f$ ,  $Join_f$ ,  $Agg_{f_m, f_a, k}$ ,  $FlatMap_f$ ,  $Sort_f$

$Map_f$ 

Transforme chaque élément de la collection via  $f$

- $f : \tau \rightarrow \tau'$
  - $Map_f : [\tau] \rightarrow [\tau']$
- 
- |   |  |
|---|--|
| <ul style="list-style-type: none"><li>• SQL : SELECT</li><li>• Mongo agg : \$project ou \$replaceWith</li></ul> | <ul style="list-style-type: none"><li>• Javascript :<br/>Array.prototype.map</li><li>• Python : map(f, ...)</li><li>• OCaml : List.map</li><li>• Java : Stream.map</li></ul> |
|---|--|
- Reformuler  $f$  avec des

Reformuler  $f$  avec des expressions définissant des attributs

$Filter_f$ 

Conserve uniquement certains éléments, choisis par  $f$

- $f : \tau \rightarrow bool$
- $Filter_f : [\tau] \rightarrow [\tau]$
- SQL : WHERE  
Reformuler  $f$  avec des conditions sur les attributs
- Mongo agg : \$match  
Reformuler  $f$  sous forme de conditions (*query*)
- Javascript :  
Array.prototype.filter
- Python : filter( $f$ , ...)
- OCaml : List.filter
- Java : Stream.filter

## $Join_f$

Combine les éléments de deux collections.  
 $e_1$  est combiné avec  $e_2$  si  $f(e_1)(e_2) = true$

- $f : \tau_1 \rightarrow \tau_2 \rightarrow bool$
- $Join_f : [\tau_1] \rightarrow [\tau_2] \rightarrow [< left : \tau_1, right : \tau_2 >]$
- SQL : JOIN  
Reformuler  $f$  sous forme de conditions dans le ON
- Mongo agg : \$lookup  
Reformuler  $f$  sous forme de conditions (*query*), et faire suivre d'un \$unwind
- Javascript, Python, OCaml, Java : pas de codage direct
- `coll1.map(e1 => coll2.filter(e2 => f(e1,e2)).map(e2 => {"left": e1, "right": e2})).flat()`



$$Agg_{f_m, f_a, k}$$

Regroupe les éléments  $x$  de la collection selon les valeurs de  $k(x)$ .

Crée un record pour chaque groupe ayant :

- la valeur de  $k(x)$  du groupe dans le champ *key*
  - la valeur de  $f_a \circ Map_{f_m}$  appliquée au groupe dans le champ *value*
- $f_m : \tau \rightarrow \tau_a$
  - $f_a : [\tau_a] \rightarrow \tau_b$
  - $k : \tau \rightarrow \tau_g$
  - $Agg_{f_m, f_a, k} : [\tau] \rightarrow \langle key : \tau_k, value : \tau_b \rangle$

## $Agg_{f_m, f_a, k}$ dans les langages concrets

- SQL : GROUP BY + fonctions d'aggrégations
- MongoDB : \$group
- Javascript : à recoder
- Python : `itertttools.groupby`
- OCaml : à recoder
- Java : `Collectors.groupingBy`

## $ReduceByKey_{f,a_k,a_v}$ : cas particulier de $Agg_{f_m,f_g,k}$

Version où l'aggrégation se fait élément par élément, deux à deux en combinant avec des résultats intermédiaires *de même type*.

- $f : \tau_v \rightarrow \tau_v \rightarrow \tau_v$
- $ReduceByKey_{f,a_k,a_v} : [\tau] \rightarrow \langle key : \tau_k, value : \tau_v \rangle$ ,  
avec  $\tau \preceq \langle a_k : \tau_k, a_v : \tau_v \rangle$

## $ReduceByKey_{f,a_k,a_v}$ : codage

$$ReduceByKey_{f,a_k,a_v} = Agg \ \lambda x.(x.a_v), \ (reduce \ (\lambda x \lambda y.(f \ x \ y))), \ \lambda x.(x.a_k)$$

avec

```
let reduce agg l =  
  let rec fold acc l' =  
    match l' with  
    | [] -> acc  
    | x::l2 -> fold (agg acc x) l2  
  in match l with  
  | [] -> Erreur  
  | x::l3 -> fold x l3
```

## *FlatMap<sub>f</sub>*

Produit pour chaque élément de la collection initiale des valeurs.  
Le résultat est la collection de toutes les valeurs produites.

- $f : \tau \rightarrow [\tau']$
- $FlatMap_f : [\tau] \rightarrow [\tau']$
- SQL : Possible avec certaines fonctions particulières (e.g. `unnest` en PostgreSQL)
- MongoDB : `$unwind`
- Javascript :  
`Array.prototype.flatMap()`
- Python :  
`itertools.chain.from_iterable`
- OCaml : `List.flatten` combiné avec `List.map`
- Java : `Stream.flatMap`

# $Sort_f$

Trie la collection selon la fonction de comparaison  $f$

- $f : \tau \rightarrow \tau \rightarrow bool$
- $Sort_f : [\tau] \rightarrow [\tau]$
- SQL : ORDER BY
- MongoDB : \$sort
- Javascript :  
    Array.prototype.sort
- Python : sorted
- OCaml : List.sort
- Java : Stream.sorted

# Union

## Assemble des collections

- $Union : [\tau] \rightarrow [\tau] \rightarrow [\tau]$
- SQL : UNION
- MongoDB : \$unionWith
- Javascript :  
    `Array.prototype.concat`
- Python : `itertools.chain`
- OCaml : `List.append`
- Java : `Stream.concat`

## Diff

## Différence entre collections

- $Diff : [\tau] \rightarrow [\tau] \rightarrow [\tau]$
- SQL : MINUS, NOT IN, NOT EXISTS
- MongoDB : parfois recodable
- Javascript, Python OCaml, Java : recoder avec *Filter* et une fonction de test d'appartenance



# Exemples

Exemples avec MongoDB aggregation pipeline

# Données à la demande : implémentation basée sur les itérateurs

- Itérateur : objet/fonction/methode fournissant les éléments un par un
- Approche naturelle pour traiter des collection lues depuis des fichiers
- Plus compliqué pour les tris, jointures, calculs de groupes :
  - on perd l'aspect flux ;
  - peut nécessiter la *matérialisation* d'une collection

# Digression : Itérateurs en Python

## Principe

- Objet avec état utilisé pour itérer sur une collection (possiblement virtuelle)
- `next()` méthode qui renvoie (ou *yields*) le prochain élément
  - throws `StopIteration` lorsqu'il n'y a plus d'éléments.

# Digression : Générateurs en Python

- “fonction” spéciale qui crée un itérateur
- `yield` statement :
  - chaque utilisation de `yield` fourni la valeur qui sera renvoyée par le prochain appel à `next()`

## Digression : Exemple de générateur

```
def foo():  
    print("begin")  
    for i in range(3):  
        print("before_yield", i)  
        yield i  
        print("after_yield", i)  
    print("end")  
  
for i in foo():  
    print("obtained", i)
```

```
begin  
before yield 0  
obtained 0  
after yield 0  
before yield 1  
obtained 1  
after yield 1  
before yield 2  
obtained 2  
after yield 2  
end
```

# Générateur pour $Map_f$

```
def op_map(f, coll):  
    for elt in coll:  
        yield f(elt)
```

## Générateur pour $\text{Filter}_f$

```
def op_filter(f, coll):  
    for elt in coll:  
        if f(elt):  
            yield elt
```

## Générateur pour $Join_f$

```
def op_join(f, coll1, coll2):  
    # Nested loops  
    for elt1 in coll1:  
        for elt2 in coll2:  
            if f(elt1, elt2):  
                yield {"left": elt1, "right": elt2}
```



# Générateur pour $Agg_{f_m, f_a, k}$

```
def op_agg(f_m, f_a, k, coll):  
    groups = dict()  
    for elt in coll:  
        key = k(elt)  
        if key not in groups:  
            groups[key] = []  
        groups[key].append(elt)  
    for k in groups:  
        yield {"key": k, "value": f_a(map(f_m, groups[k]))}
```

# Générateur pour $Agg_{f_m, f_a, k}$ , utilisant le tri

```
def op_agg(f_m, f_a, k, coll):  
    coll = sorted(coll, k)  
    key = None  
    grp = []  
    for elt in coll:  
        key_elt = k(elt)  
        if key != key_elt and key is not None:  
            yield {"key": key, "value": f_a(map(f_m, grp))}  
            key = key_elt  
            grp = []  
        elif key is None:  
            key = key_elt  
        grp.append(elt)  
    if key is not None:  
        yield {"key": key, "value": f_a(map(f_m, grp))}
```

## Générateur pour $\text{FlatMap}_f$

```
def op_flatmap(f, coll):  
    for elt in coll:  
        for elt2 in f(elt):  
            yield elt2
```

## Et en distribué ?

- Collection répartie sur des serveurs  $S_1, \dots, S_n$
- Résultat d'un calcul : union des résultats produits par chaque serveur
- Rien de particulier à faire pour  $Map_f, Filter_f, FlatMap_f$

## $Sort_f$ en distribué

- Trier en local sur chaque  $S_i$
- Un des serveurs est élu pour fournir le résultat (arbitrairement  $S_1$ )
- $S_1$  demande à tous les serveurs leur premier élément
- On itère ensuite, jusqu'à ce qu'il n'y ait plus de valeur :
  - déterminer  $i$  comme le numéro du serveur ayant produit la plus petite valeur
  - yield la valeur
  - mettre à jour la valeur du serveur  $i$

## $Agg_{f_m, f_a, k}$ en distribué

- Chaque valeur de  $k(x)$  se voit attribuer un serveur
  - e.g. via un hash entre 1 et  $n$
- Chaque serveur redistribut ses éléments  $x$  en les envoyant vers le serveur en fonction de la valeur de  $k(x)$
- On applique ensuite l'algorithme local

## *Join<sub>f</sub>* en distribué

Si  $f = \lambda e_1. \lambda e_2. e_1[a] = e_2[a]$

- Chaque valeur de  $a$  se voit attribuer un serveur
  - e.g. via un hash entre 1 et  $n$
- Chaque serveur redistribut ses éléments en les envoyant vers le serveur en fonction de la valeur de  $a$
- Calcule pour chaque valeur de  $a$  le sous-produit cartésien des deux sous-collections correspondantes