

# **Cours de Business Intelligence**

**Master : Traitement intelligent des systèmes**

**Préparé par:** Mme. ELANSARI Khawla

**Année Universitaire:** 2021/ 2022

# Data Warehouse (2)

1. Exercice 2
2. Rappels
3. Types de dimensions
4. Méta-Données
5. Stockage des données
6. Architecture du DW: Ralph Kimball
7. Architecture du DW: Bill Inmon

## | **Exercice 2**

1- Concevoir un modèle en étoile qui permet d'analyser les ventes d'une entreprise de restauration rapide. Le principe est de mesurer les ventes grâce aux quantités vendues et aux bénéfices, en fonction des ventes réalisées par jour, dans un restaurant donné, pour un aliment donné.

L'objectif est de pouvoir analyser les ventes par jour, par semaine, par mois et par année. Les restaurants peuvent être regroupés en fonction de leur ville et de leur pays.

## | **Exercice 2**

1- Concevoir un modèle en étoile qui permet d'analyser les **ventes** d'une entreprise de restauration rapide. Le principe est de mesurer les ventes grâce aux quantités vendues et aux bénéfices, en fonction des ventes réalisées par jour, dans un restaurant donné, pour un aliment donné.

L'objectif est de pouvoir analyser les ventes par jour, par semaine, par mois et par année. Les restaurants peuvent être regroupés en fonction de leur ville et de leur pays.

## | **Exercice 2**

1- Concevoir un modèle en étoile qui permet d'analyser les **ventes** d'une entreprise de restauration rapide. Le principe est de mesurer les ventes grâce aux **quantités vendues** et aux **bénéfices**, en fonction des ventes réalisées par jour, dans un restaurant donné, pour un aliment donné.

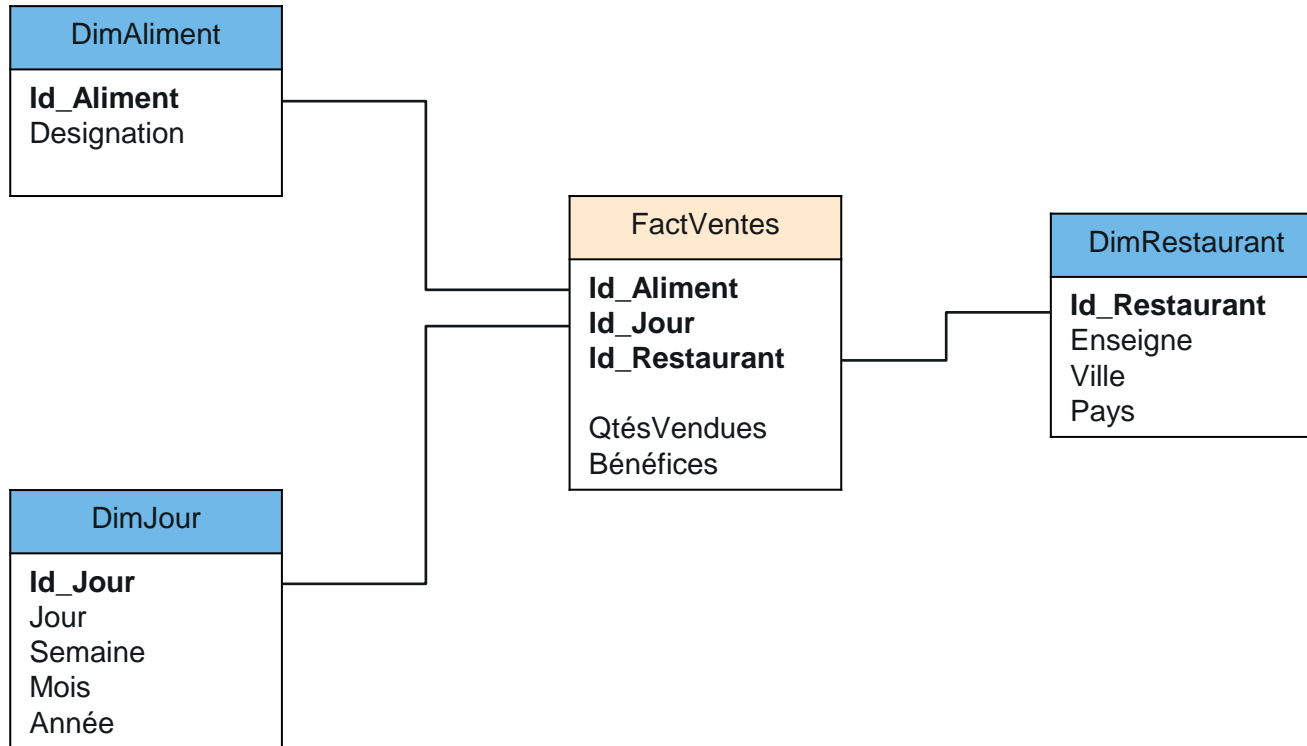
L'objectif est de pouvoir analyser les ventes par jour, par semaine, par mois et par année. Les restaurants peuvent être regroupés en fonction de leur ville et de leur pays.

## | **Exercice 2**

1- Concevoir un modèle en étoile qui permet d'analyser les **ventes** d'une entreprise de restauration rapide. Le principe est de mesurer les ventes grâce aux **quantités vendues** et aux **bénéfices**, en fonction des ventes réalisées par **jour**, dans un **restaurant** donné, pour un **aliment** donné.

L'objectif est de pouvoir analyser les ventes par jour, par semaine, par mois et par année. Les restaurants peuvent être regroupés en fonction de leur ville et de leur pays.

## | *Exercice 2*



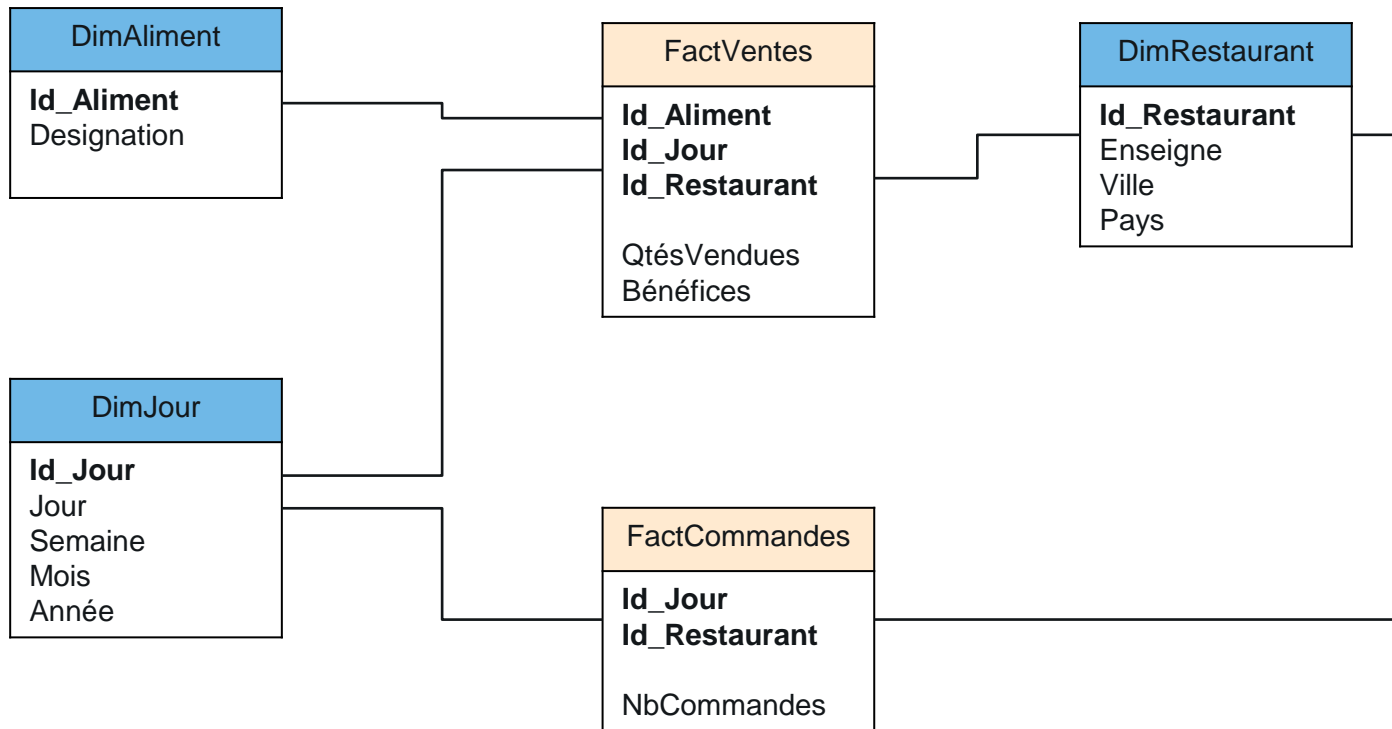
## | ***Exercice 2***

2- On souhaite à présent mesurer le nombre de commandes qui est donné par jour et par restaurant.

Etendre le modèle précédent afin de prendre en compte cet aspect.

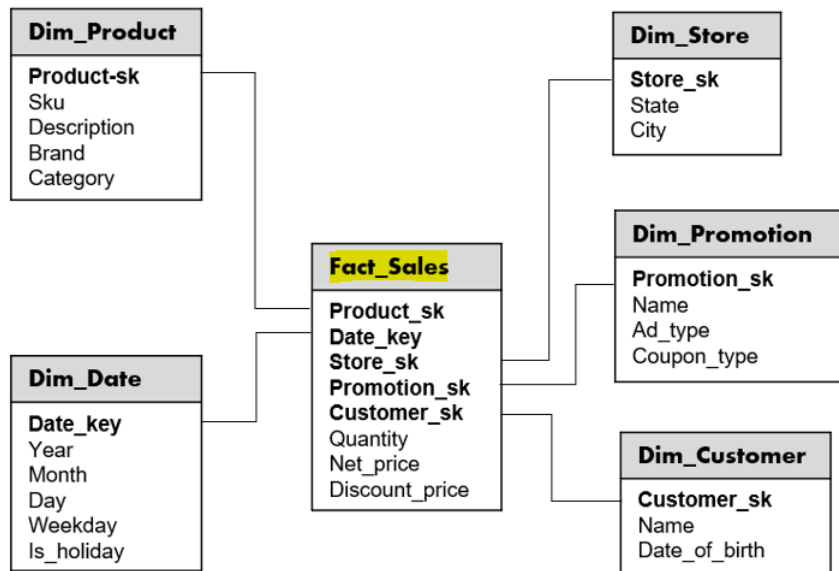


## | *Exercice 2*



## 2. Rappel : Modèle dimensionnel

### Modèle en étoile



dim\_product table

product_sk	sku	description	brand	category
30	OK4012	Bananas	Freshmax	Fresh fruit
31	KA9511	Fish food	Aquatech	Pet supplies
32	AB1234	Croissant	Dealicious	Bakery

dim\_store table

store_sk	state	city
1	WA	Seattle
2	CA	San Francisco
3	CA	Palo Alto

fact\_sales table

date_key	product_sk	store_sk	promotion_sk	customer_sk	quantity	net_price	discount_price
140102	31	3	NULL	NULL	1	2.49	2.49
140102	69	5	19	NULL	3	14.99	9.99
140102	74	3	23	191	1	4.49	3.89
140102	33	8	NULL	235	4	0.99	0.99

dim\_date table

date_key	year	month	day	weekday	is_holiday
140101	2014	jan	1	wed	yes
140102	2014	jan	2	thu	no
140103	2014	jan	3	fri	no

dim\_customer table

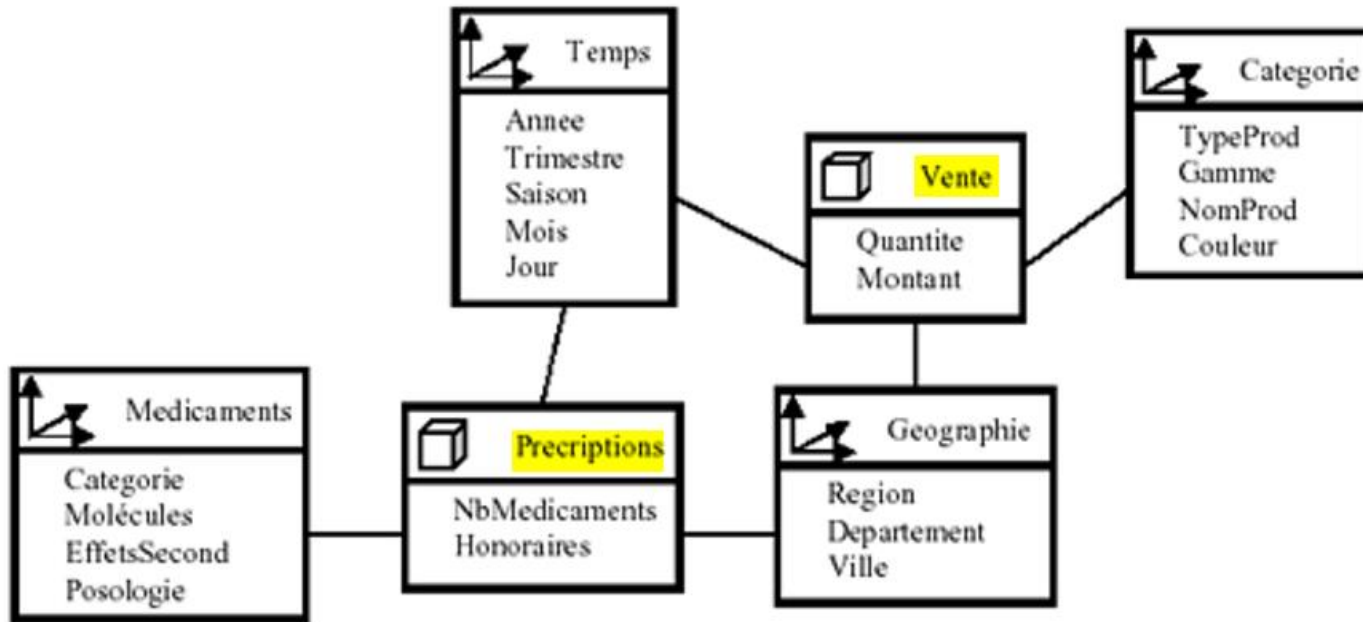
customer_sk	name	date_of_birth
190	Alice	1979-03-29
191	Bob	1961-09-02
192	Cecil	1991-12-13

dim\_promotion table

promotion_sk	name	ad_type	coupon_type
18	New Year sale	Poster	NULL
19	Aquarium deal	Direct mail	Leaflet
20	Coffee & cake bundle	In-store sign	NULL

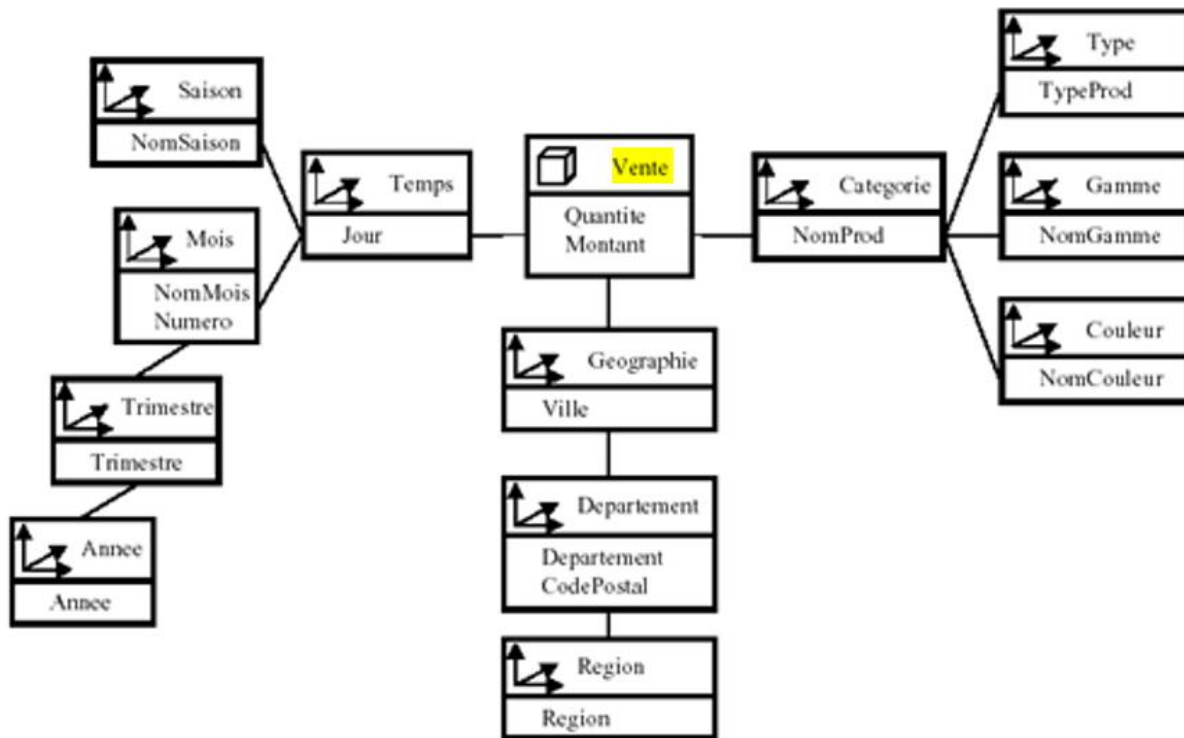
## | 2. Rappel : Modèle dimensionnel

### Modèle en constellation



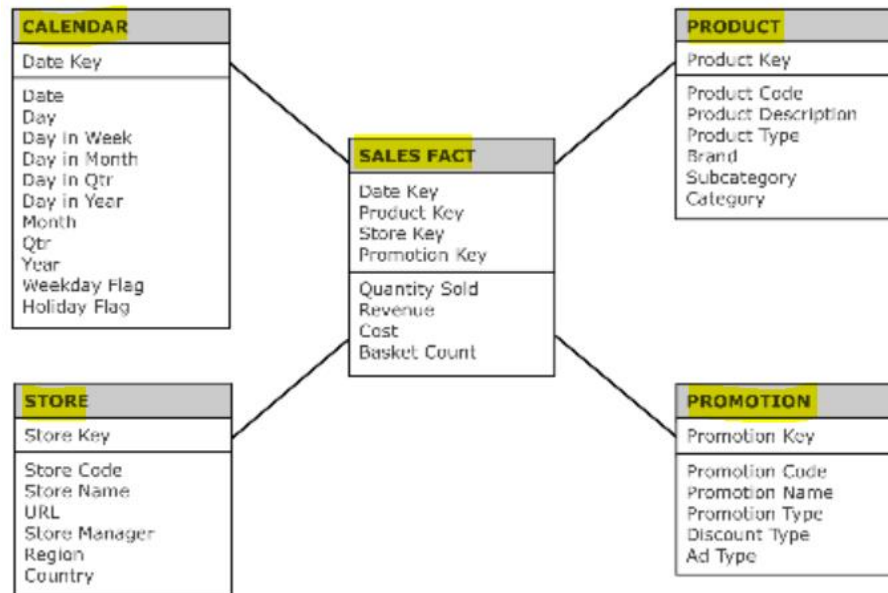
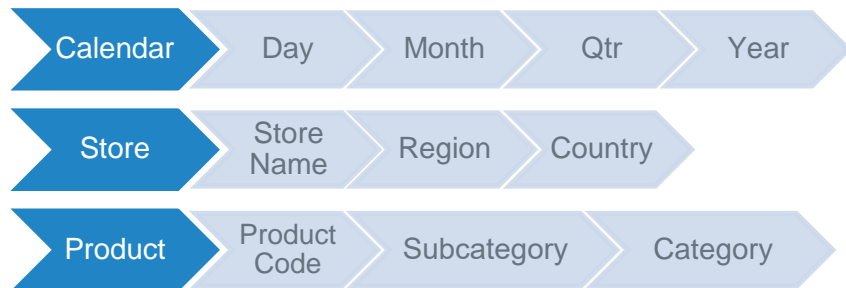
## | 2. Rappel : Modèle dimensionnel

### Modèle en Flocon



## | 2. Rappel : Granularité

La granularité est ce qui permet de définir le niveau de détail des informations présentes dans une ligne d'une table de faits. Il est défini par un ensemble minimal de dimensions.

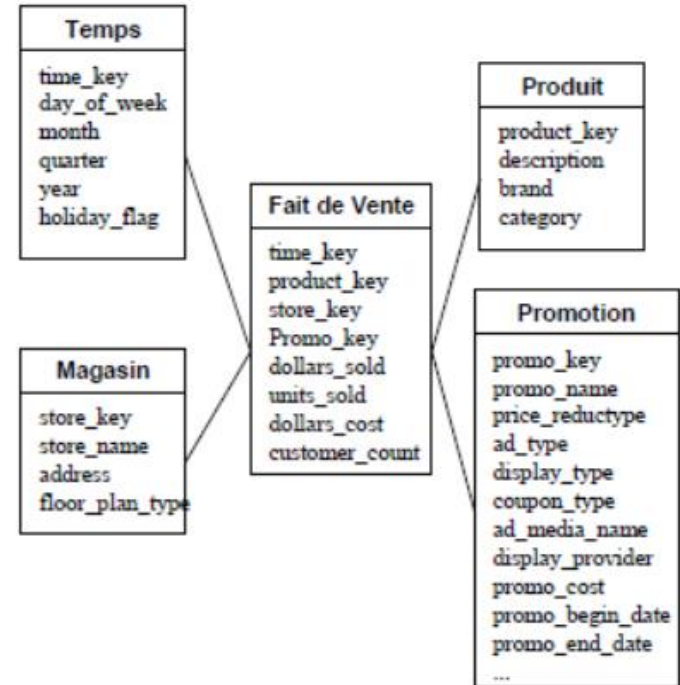


### | 3. Types de dimensions

**Dimension Conforme:** On parle de dimension conforme ou partagée lorsque la dimension est utilisée par les faits de plus qu'un data Mart. **Ex:** la dimension « Produit » est utilisée par différents data Mart « Finance », « Marketing »...

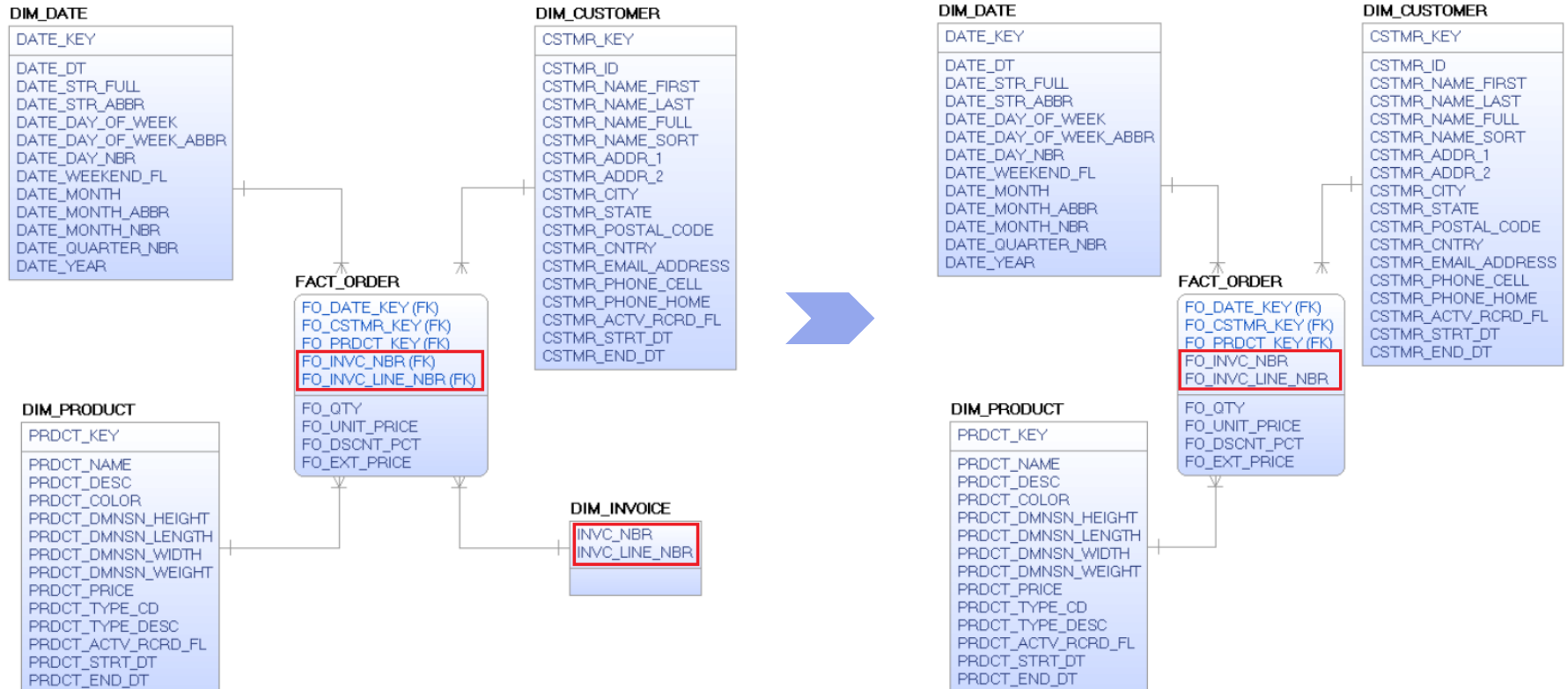
**Causale:** Dimension qui provoque le fait.

**Ex:** la dimension Promotion est supposée avoir provoqué le Fait de Vente



### 3. Types de dimensions

**Dimension Dégénérée:** Dimension sans attribut (Pas de table, mais la clé de dimension est dans la table de fait)



### | 3. Types de dimensions

**Dimension Junk:** La dimension de genre « Junk dimension » est une dimension qui contient toutes sortes de flags, statuts, codes qui ne font partie d'aucune dimension régulière.

**Ex:** Un exemple de ceci peut être la couleur de la voiture (rouge, noir, bleu, etc.) et le style de la carrosserie (Sedan, VAN, SUV, etc.).

Le nombre des possibilités est limité et, si on crée deux dimensions Couleur et Carrosserie, elles seront limitées à un seul attribut. Afin d'éliminer ces petites dimensions, nous créons une seule dimension «Junk» qui fait un cross join de tous les attributs possibles en une seule dimension qui sera utilisée dans la table des faits.

DIM_CAR_ATTRIBUTES	
P *	DCA_KEY DCA_COLOR DCA_BODY_STYLE

DCA_KEY	DCA_COLOR	DCA_BODY_TYPE
1	BLACK	SEDAN
2	WHITE	SEDAN
3	RED	SEDAN
4	SILVER	SEDAN
	...	...
101	BLACK	SUV
102	WHITE	SUV
103	RED	SUV
104	SILVER	SUV
	...	...



### | 3. Types de dimensions

**Dimension à évolution lente:** Il s'agit de dimensions évoluant peu au regard de l'application : changement de nom d'un client, changement de dénomination d'une région géographique... Ce type d'évolution est la plus courante.

**Ex:**

- Un client peut changer d'adresse, se marier (et donc changer de nom) ;
- Un produit peut changer de nom, de formulation (« Yaourt à la vanille » devient « saveur Vanille ») ;
- Une entreprise peut changer son organisation commerciale (secteurs) ou sa nomenclature de produits.

Il est possible de gérer cette situation en choisissant entre ces solutions :

- **Sol 1** : Écrasement de l'ancienne valeur / remplacer
- **Sol 2** : Création d'une nouvelle ligne
- **Sol 3** : Valeur d'origine / valeur courante

### | 3. Types de dimensions (SCD)

**Sol 1:** Écraser l'ancienne valeur / La remplacer.

Cette solution est la plus simple à mettre en œuvre, car elle ne nécessite ni modélisation particulière, ni besoin d'ajouter de nouveaux enregistrements dans la table. Elle consiste à mettre à jour les données, en écrasant l'ancienne version

Id	EAN Code	Product Name	Brand	Product Category
1	977147396801	Canon EOS Rebel	Cannon	Camera
2	977147396802	Nikon Coolpixx	Nikon	Camera
3	977147396803	Sony Cyber-shot	Sony	Camera
4	977147396804	Olympus XZ-1	Olympus	Camera

Id	EAN Code	Product Name	Brand	Product Category
1	977147396801	Canon EOS Rebel	Cannon	Camera
2	977147396802	Nikon Coolpix	Nikon	Camera
3	977147396803	Sony Cyber-shot	Sony	Camera
4	977147396804	Olympus XZ-1	Olympus	Camera

### | 3. Types de dimensions (SCD)

**Sol 2:** Ajouter une ligne avec un nouvel ID pour la nouvelle valeur. Aussi, on peut ajouter deux colonnes date de début et de fin d'utilisation à la table pour garder une trace sur les dates des changements ou ajouter un flag (Y, N) pour spécifier l'enregistrement 'Actif'. => Utilisées pour garder tout l'historique des données.

N.B: Une nouvelle ligne est insérée à chaque fois un changement a lieu sur les tables source

Id	EAN_Code	Product_Name	Brand	Product_Category
1	977147396801	Canon EOS Rebel	Cannon	Camera
2	977147396802	Nikon Coolpix	Nikon	Camera
3	977147396803	Sony Cyber-shot	Sony	Camera
4	977147396804	Olympus XZ-1	Olympus	Camera
5	977147396804	Olympus XZ-1	Olympus	Electronics

### | 3. Types de dimensions (SCD)

**Sol 3:** Ajouter une colonne à la table de dimension pour chaque colonne dont on souhaite garder une trace des changements, afin de stocker l'historique des modifications. => L'ancienne valeur n'est utile que pendant un certain temps pour étudier les effets d'une transition.

Id	EAN Code	Product Name	Brand	Cat Current	Cat Previous
1	977147396801	Canon EOS Rebel	Cannon	Camera	Camera
2	977147396802	Nikon Coolpix	Nikon	Camera	Camera
3	977147396803	Sony Cyber-shot	Sony	Camera	Camera
4	977147396804	Olympus XZ-1	Olympus	Electronics	Camera

### | 3. Types de dimensions (RCD)

**Dimension à évolution rapide:** Il s'agit des dimensions possédant des attributs variant fréquemment : salaire ou poste des collaborateurs...

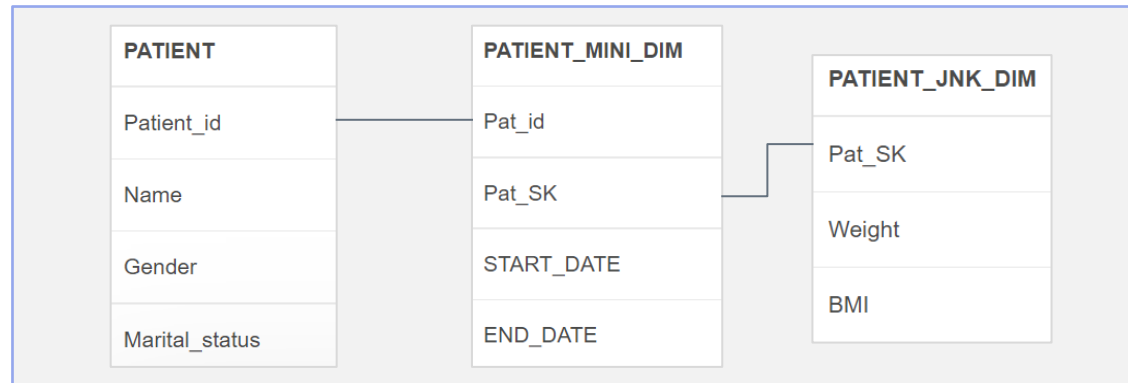
**Ex:** Considérons la dimension **PATIENT**. Les attributs : Patient\_id, Name, Gender et Marital\_Status changent rarement, cependant les attributs Weight et BMI changent fréquemment (à chaque consultation) => La solution serait de séparer les attributs qui changent fréquemment des attributs qui changent rarement.

PATIENT
Patient_id
Name
Gender
Marital_status
Weight
BMI



### | 3. Types de dimensions (RCD)

PATIENT
Patient_id
Name
Gender
Marital_status
Weight
BMI



## | 4. Méta-Données

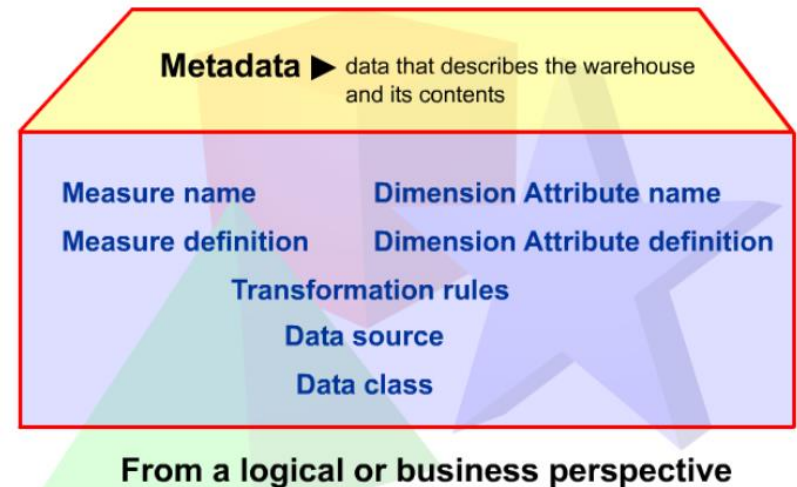
Essentielles pour la gestion des couches de données et des processus de construction

- Aide à l'administrateur et au concepteur

- architecture complexe
- gros volumes de données
- nombreux processus évolutifs au cours du temps

- Types de méta-données

- **data dictionary** : définitions des schémas des BD
- **rafraîchissement**: structure et fréquence de l'alimentation
- **transformations** : définition et flux
- **versions** : contrôle des changements de méta-données
- **statistiques** : profils des données entreposées
- **sécurité** : conditions d'accès aux données
- **localisation physique des données**



## | 5. *Stockage des données*

Le stockage au sein d'un datawarehouse a un besoin de synthèse (agrégation des données) et un besoin de détails (conservation des données détaillées).

Ce stockage peut être réalisé de trois manières différentes :

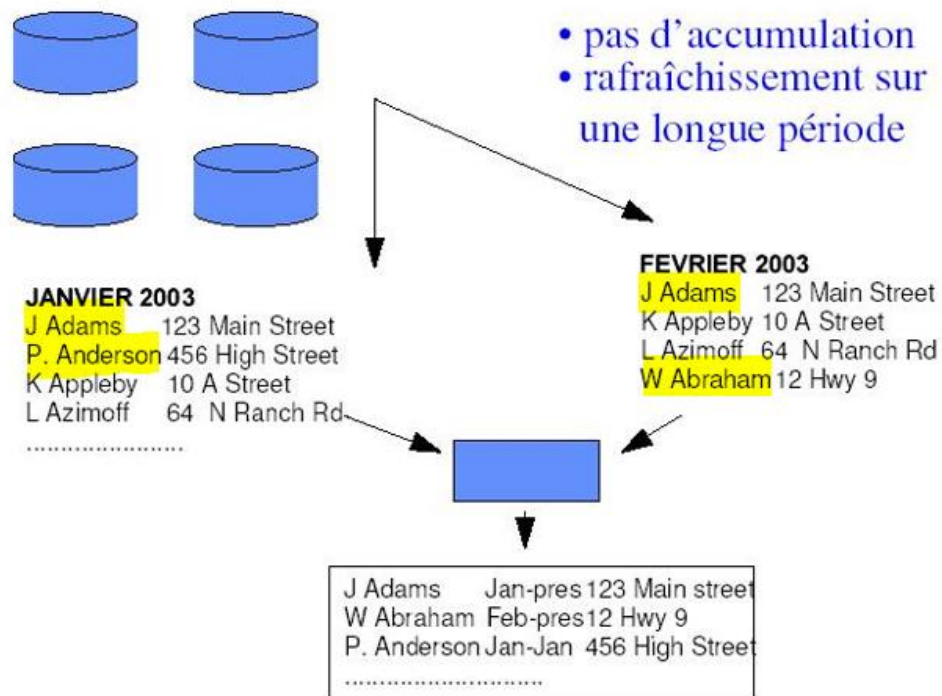
- structure directe simple,
- structure de cumul simple,
- par résumé déroulant.



## | 5. Stockage des données

### *Structure directe simple*

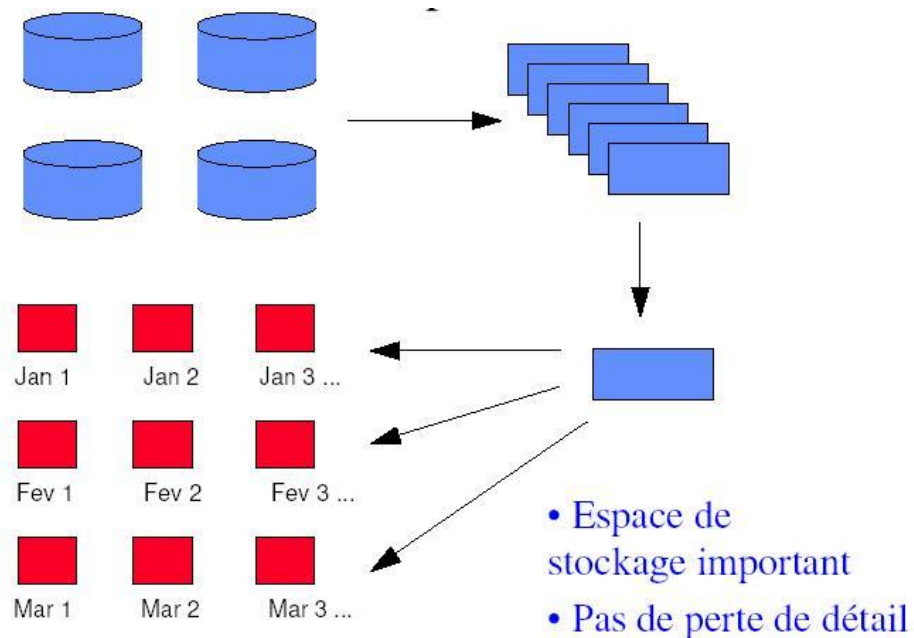
On fait des mises à jour du datawarehouse avec des laps de temps important.



## | 5. Stockage des données

### *Structure de cumul simple*

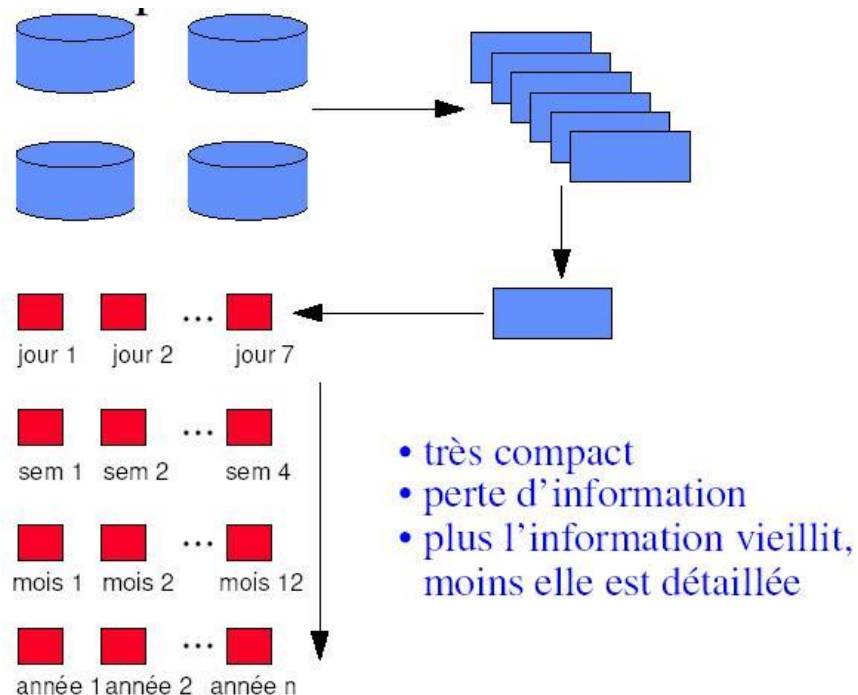
On stocke les données de chaque mise à jour, les mises à jour étant fréquentes (par exemple tous les jours) on a un espace occupé important, mais on ne perd pas d'information.



## | 5. Stockage des données

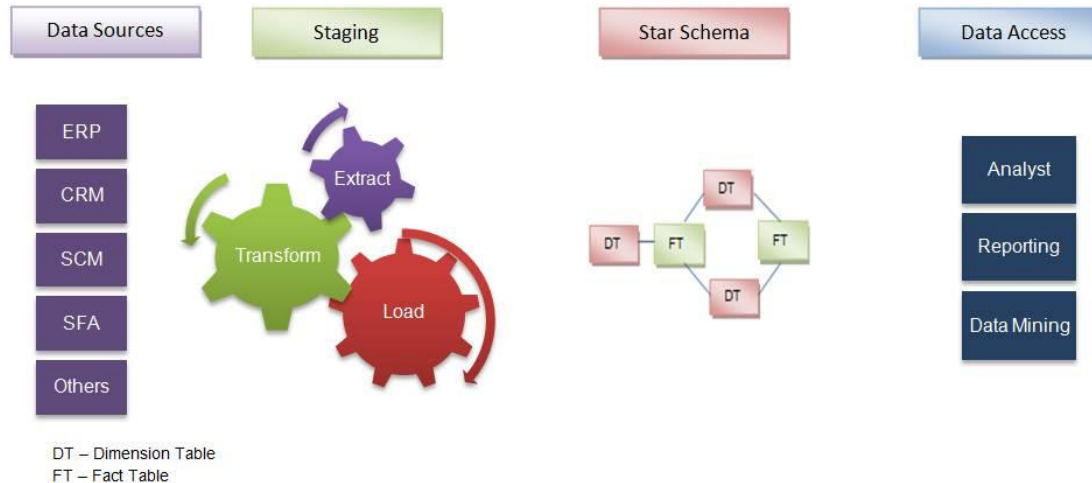
### *Structure par résumé déroulant*

A chaque mise à jour, on stocke des données détaillées, et on synthétise les anciennes données en fonction de leur âge. Plus une donnée est vieille, moins elle est détaillée

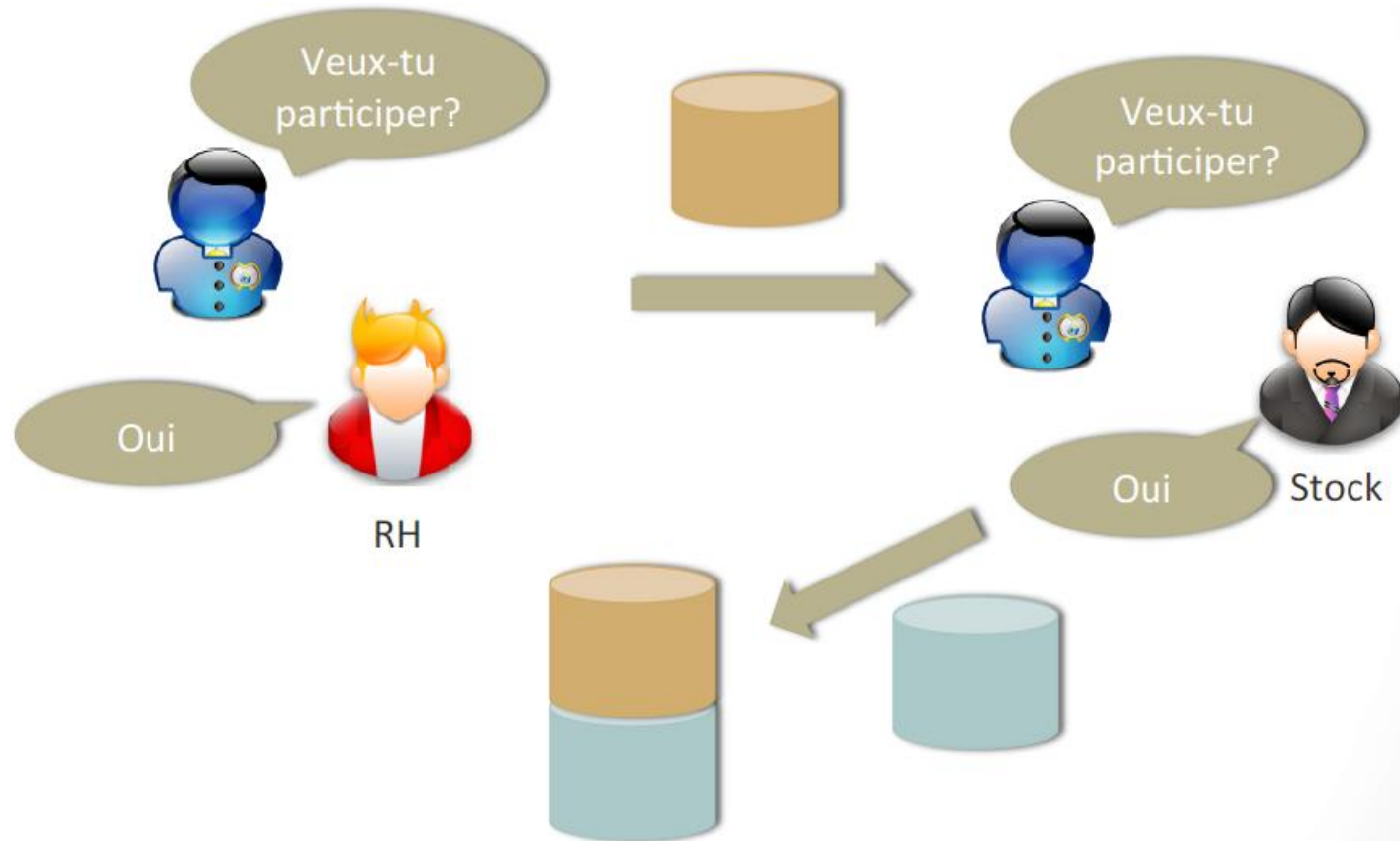


## | 6. Architecture du DW: Ralph Kimball

- Cette approche commence par la reconnaissance des processus métier et des questions auxquelles le Datawarehouse doit répondre.
- Documenter toutes les sources de données disponibles
- Créer des pipelines ETL qui extraient, transforment et chargent les données des sources de données dans un modèle de données dénormalisé. Le modèle dimensionnel est construit sous la forme d'un schéma en étoile ou d'un schéma en flocon
- Le modèle dimensionnel est généralement construit autour et au sein de Data Marts pour les différents départements.

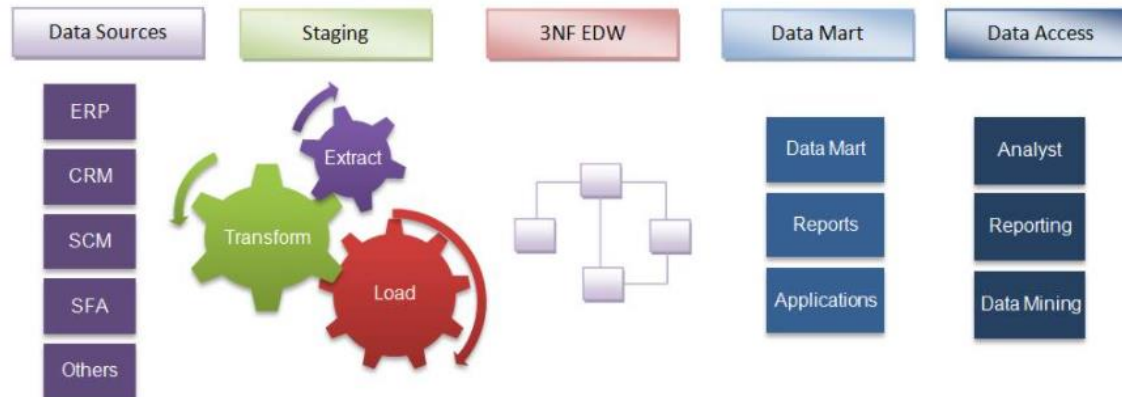


## | 6. Architecture du DW: Ralph Kimball

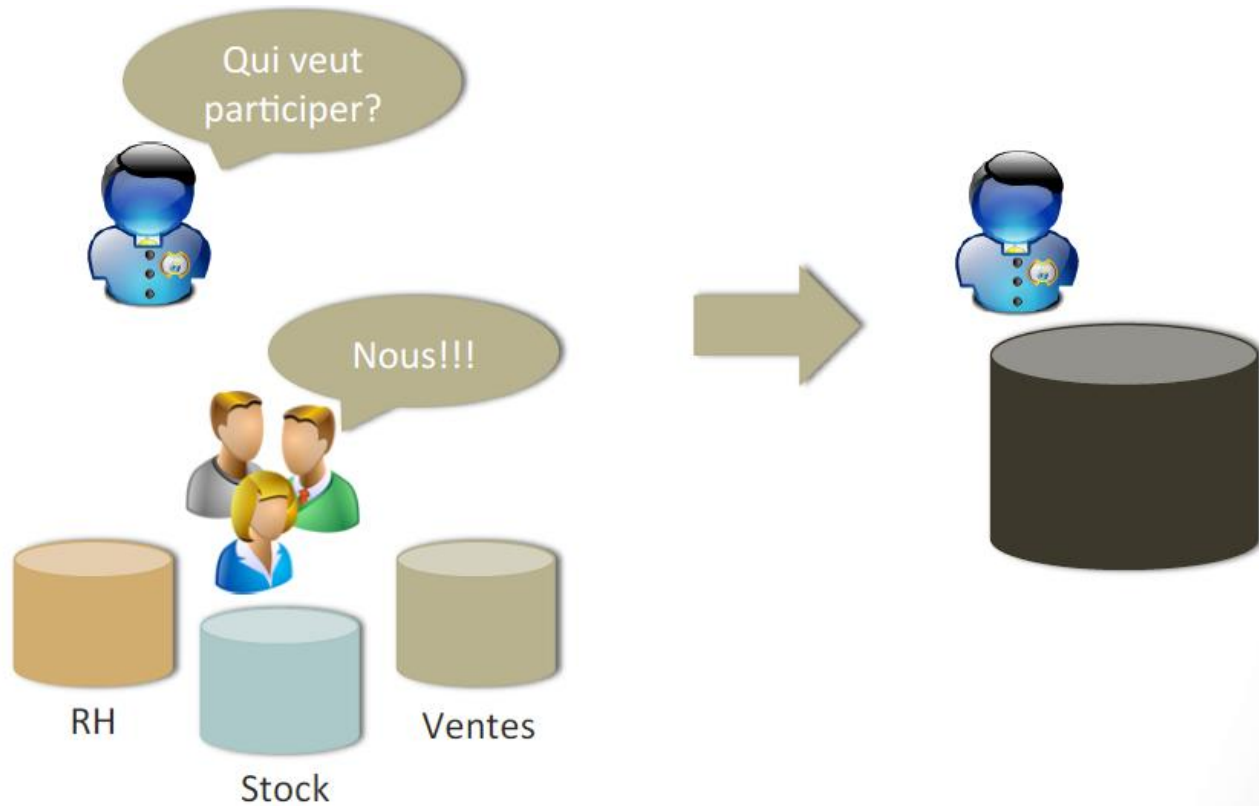


## | 7. Architecture du DW: Bill Inmon

- A partir du modèle logique, construire le modèle physique. Utiliser des processus ETL pour extraire des données de différentes sources, transformer des données et les charger dans un modèle de données normalisé. Le modèle de données normalisé est le cœur de l'entrepôt de données.
- Créer des Data Marts pour les différents départements. Les données sont accessibles via des Data Marts pour tous les besoins de reporting, l'entrepôt de données agit comme une source unique de vérité.



## | 7. Architecture du DW: Bill Inmon



# Analyse Multidimensionnelle

1. Analyse Multidimensionnelle
2. Cube OLAP
3. Les opérations multidimensionnelles
4. Implémentation et déploiement d'un cube OLAP
5. Extraction des connaissances
6. Visualisation des données



## | 1. *Analyse Multidimensionnelle*

Permet d'obtenir des informations déjà agrégées selon les besoins de l'utilisateur : simplicité et rapidité d'accès.

Les données sont intégrées selon différentes dimensions

**Ex:** analyse des ventes /catégorie de produit 1D.

+ /année 2D.

+ /département commercial 3D.

+ / zone géographique 4D.

## | 2. *Cube OLAP*

Un **cube OLAP** est une **structure de données multidimensionnelle** stockant les faits comme des mesures indexées par plusieurs dimensions. Ainsi, chaque cellule d'un cube représente la mesure ou valeur quantitative d'un fait sur le croisement de plusieurs dimensions.

La mesure peut être additive, semi-additive ou non additive :

**Mesure additive** - on peut sommer sur toutes les dimensions tout en conservant un sens.

Ex: si l'on considère les sommes des ventes par produits, villes et mois, on peut faire la somme des valeurs des cellules tout en conservant un sens aux données.

**Mesure non additive** - on ne peut pas sommer les valeurs des cellules en conservant un sens.

Ex: si le cube contient des moyennes de ventes par mois, produit et ville, il n'y a pas de sens à sommer les valeurs des cellules des villes pour les regrouper en départements.

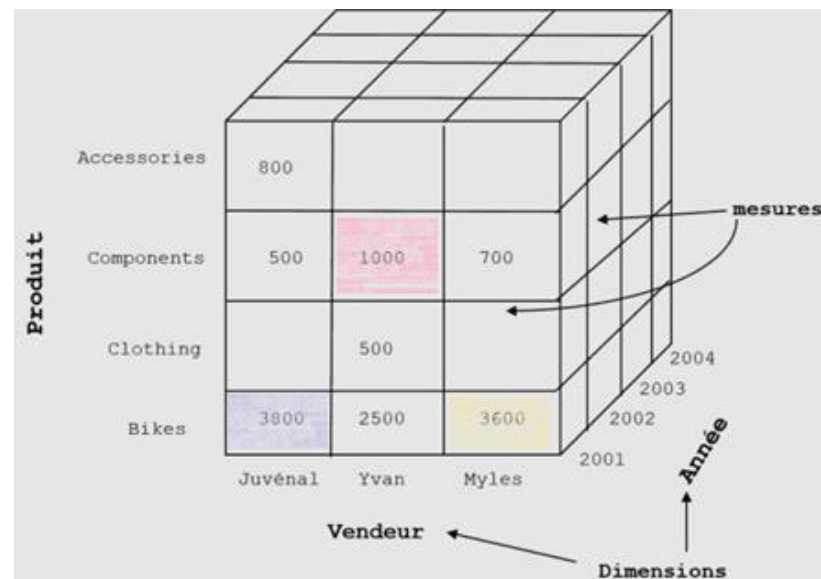
**Mesure semi-additive** - on peut sommer sur certaines dimensions en gardant un sens mais pas sur toutes.

Ex: si on considère un cube décrivant l'état des stocks par ville, produit et mois, il est possible de faire la somme sur les dimensions ville et produit pour connaître l'état global des stocks pour toutes les villes ou pour tous les produits, mais il n'y a aucun sens à sommer sur la dimension temporelle des mois.

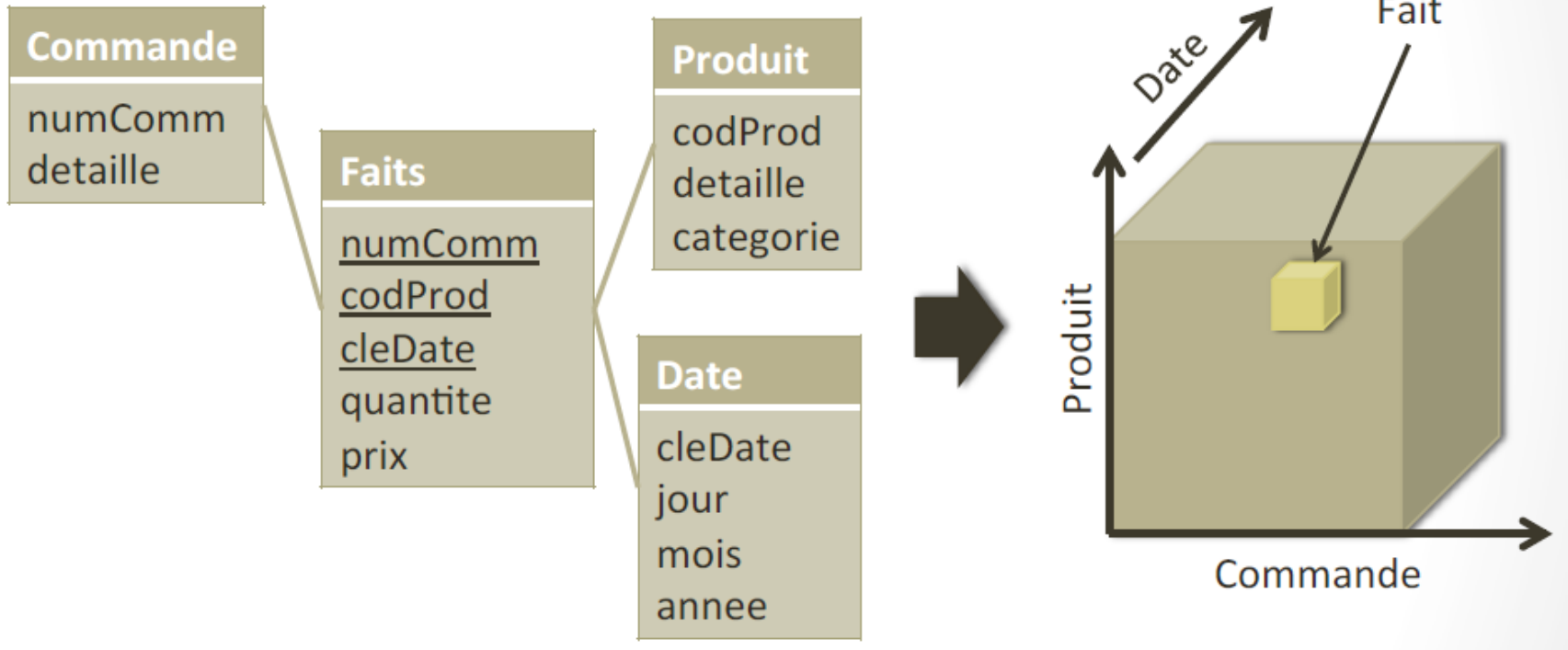
## | 2. Cube OLAP

Supposons que nous souhaitons calculer la **somme des ventes** par vendeur, produit et par année. La représentation de ce tableau de données sous forme d'une structure multidimensionnelle fournit le cube OLAP suivant.

Vendeur	Produit	Date de vente	Prix de vente
Juvénal	Accessories	01/04/2001	800
Myles	Bikes	09/05/2001	1400
yvan	Clothing	02/02/2002	500
yvan	Components	02/03/2002	1000
Juvénal	Bikes	15/03/2002	1800
Juvénal	Bikes	10/03/2003	2000
Myles	Components	12/10/2003	700
Myles	Bikes	25/12/2003	2200
Juvénal	Components	10/01/2004	500
...	...	...	...
yvan	Bikes	15/11/2004	2500



## | 2. Cube OLAP

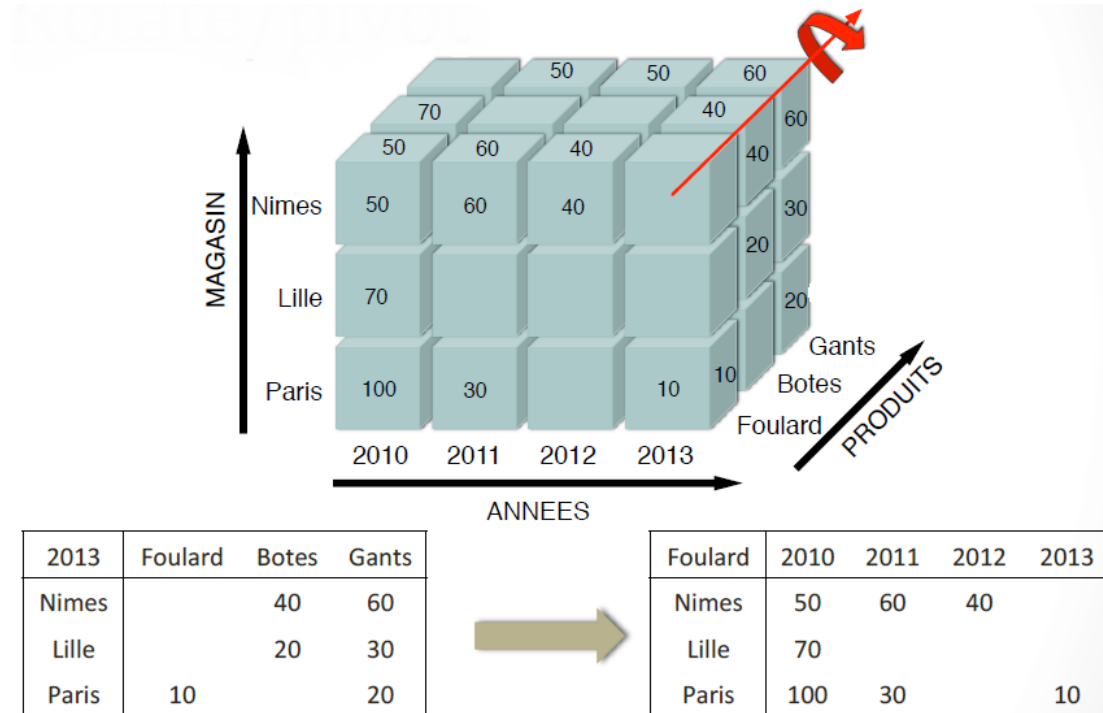


### | 3. Les opérations multidimensionnelles

#### *Manipulation d'un cube OLAP*

Les cubes multidimensionnels disposent de 3 opérateurs multidimensionnels pour leur exploitation:

**La rotation du cube (ROTATE/SWITCH)** - Rotation à 90° de 2 dimensions du cube.

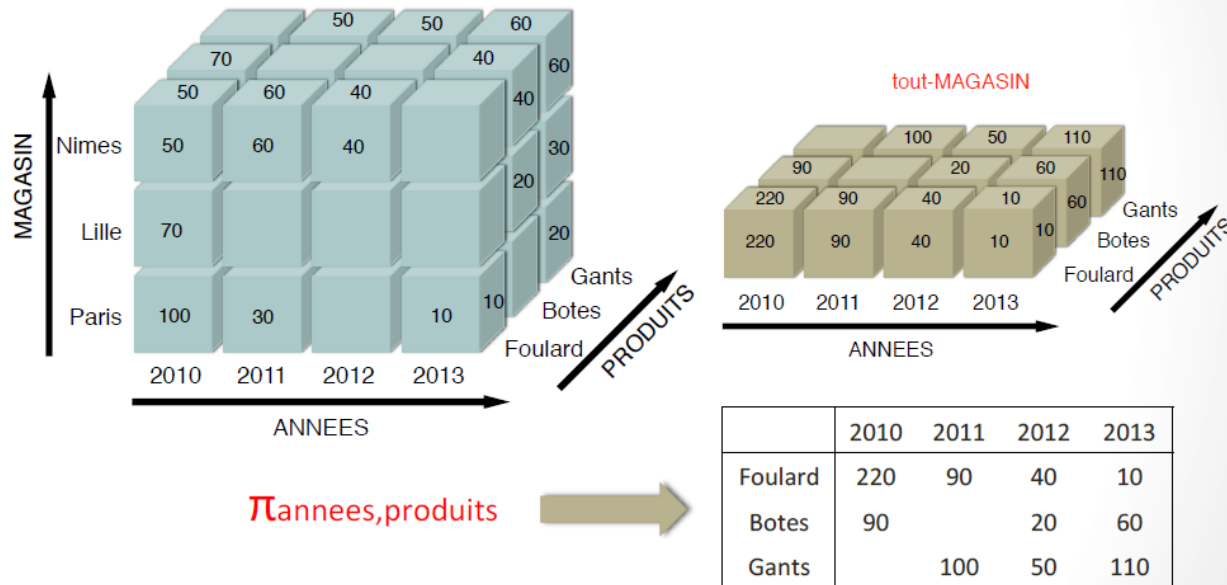


### | 3. Les opérations multidimensionnelles

#### Manipulation d'un cube OLAP

**Extraction du cube (SLICING/DICING) :** Extraire du cube un bloc de données correspondant à un croisement entre plusieurs dimensions (Sélection sur les mesures/ Les dimensions).

**SLICING** - Projection selon une dimension du cube

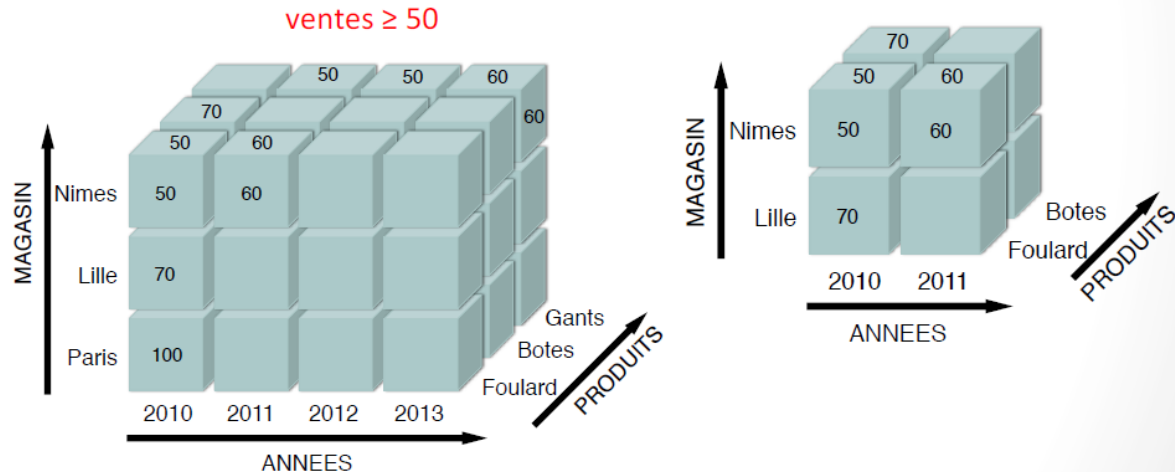


### | 3. Les opérations multidimensionnelles

#### *Manipulation d'un cube OLAP*

**DICING** - consiste à extraire un bloc de mesures en s'appuyant sur des critères d'attributs de dimensions.

Ex: Extraire les données des années 2010 et 2011 dans les magasins de Lille et Nimes en gardant tous les attributs de la dimension Produit excepte l'attribut « Gants »



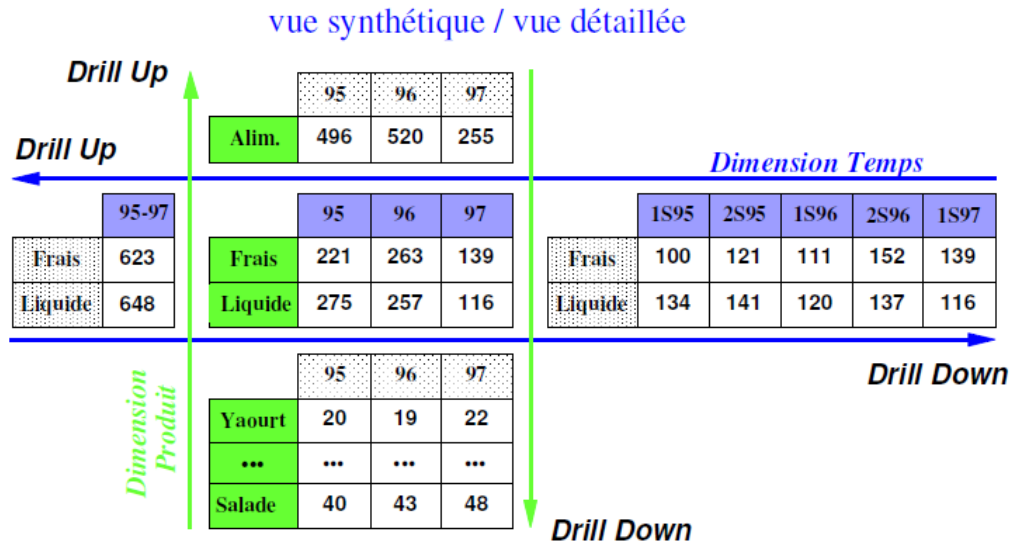
### | 3. Les opérations multidimensionnelles

#### Manipulation d'un cube OLAP

##### ROLL UP/ DRILL DOWN

**ROLL UP** – représente les données du cube à un niveau de granularité supérieur conformément à la hiérarchie définie sur la Dimension

**DRILL DOWN** – représente les données du cube à un niveau de granularité de niveau inférieur, donc sous une forme plus détaillée





## | 4. *Implémentation et déploiement d'un cube OLAP*

### *Stockage physique*

Les technologies qui sont capables de stocker les données de façon multidimensionnelle et répondre aux requêtes multidimensionnelles sont couramment appelées les **technologies OLAP**.

Le principal problème posé par les structures de données multidimensionnelles est leur nature **peu dense**, et **très éparse**. De très nombreuses cellules sont vides, les données ne sont pas distribuées uniformément dans tout l'espace multidimensionnel, elles sont concentrées en groupes dans des espace-temps où les événements métier occurrent le plus souvent (Ex: en période de soldes les ventes augmentent). Ce problème rend le stockage et le traitement des données du cube plus complexe.

Il existe trois stratégies de stockage physique : le stockage sous la forme relationnelle (**ROLAP**), sous la forme multidimensionnelle (**MOLAP**) ou une solution hybride (**HOLAP**) combinant ces deux premières approches.

## | 4. *Implémentation et déploiement d'un cube OLAP*

### *Stockage physique*

#### ROLAP

On nomme ROLAP l'**approche Relationnel OLAP**. Les données sont stockées sous la forme de tables relationnelles. Elles sont **modélisées sous la forme de schémas en étoile ou flocon**. Les requêtes multidimensionnelles doivent alors être traduites en requêtes relationnelles (SQL).

Ce modèle est excellent vis à vis de la capacité de stockage, mais les requêtes sont difficiles à définir et à mettre en œuvre et sont coûteuses.

L'outil **SSAS (SQL Server Analysis Services)** de Microsoft est un exemple de logiciel qui implémente l'approche ROLAP

## | 4. *Implémentation et déploiement d'un cube OLAP*

### *Stockage physique*

#### **MOLAP**

On nomme MOLAP l'**approche Multidimensionnelle OLAP**. La technologie de stockage est multidimensionnelle. Les données sont stockées sous la forme de tableaux multidimensionnels, des index multidimensionnels sont définis. Elle utilise des techniques de compression face à la faible densité des données. La taille des données pouvant être ainsi stockées est faible par rapport à la solution ROLAP. Cependant, les requêtes sont écrites de manière intuitive et efficace.

- On trouve en colonne tous les axes, puis tous les indicateurs
- Chaque cellule du cube est stockée par une ligne dans la matrice

Le langage utilisé pour la manipulation de ces données est le **MDX (Multi Dimensional eXpression)**. Dans le domaine du contrôle de gestion et en finance par exemple, ce sont les moteurs MOLAP qui sont utilisés ; Hyperion ESSBASE, racheté par Oracle, est un exemple de solution MOLAP ;

## | 4. *Implémentation et déploiement d'un cube OLAP*

### *Stockage physique*

#### HOLAP

On nomme HOLAP l'**approche Hybride OLAP**. Cette technologie combine les deux solutions précédentes. Les données sont stockées dans une base de données relationnelle, et les calculs sont faits dans une base multidimensionnelle.

## | 5. *Extraction des connaissances*

### *OLAP vs. Data Mining*

**Data Mining** : l'utilisateur cherche des corrélations non évidentes

Ex: «Quelles sont les caractéristiques de l'achat de yaourts ?»

**OLAP** : l'utilisateur cherche à confirmer des intuitions

Ex: «A-t-on vendu plus de yaourts en Région Parisienne qu'en Bretagne en 2003 ?»

#### Data Mining

Multidimensional Analysis.

It has large number of dimensions.

Deals with the summary of data.

Insight and Prediction.

It is used to predict the future.

Bottom-up approach.

Discovery driven.

It is a emerging technique.

#### OLAP

Online analytical processing.

It has limited number of dimensions.

Deals with the detailed Transaction level data.

Analysis.

It is used to analyze the past.

Top-down approach.

Query driven.

Widely Used.

## | 5. *Extraction des connaissances*

### *OLAP vs. Data Mining*

Il existe 2 types de techniques d'extraction de données:

#### Techniques à base d'**OLAP**

- Requêteurs - donne une réponse à une question plus ou moins complexe (type SQL)
- EIS (Executive Information Systems) outils de visualisation et de navigation dans les données statistiques + interfaçage graphique
- Applications spécialisées (ad-hoc) applications développées spécialement pour les besoins de l'entreprise

#### Techniques à base de **Data Mining**

- outils évolués de prédiction, simulation, ...

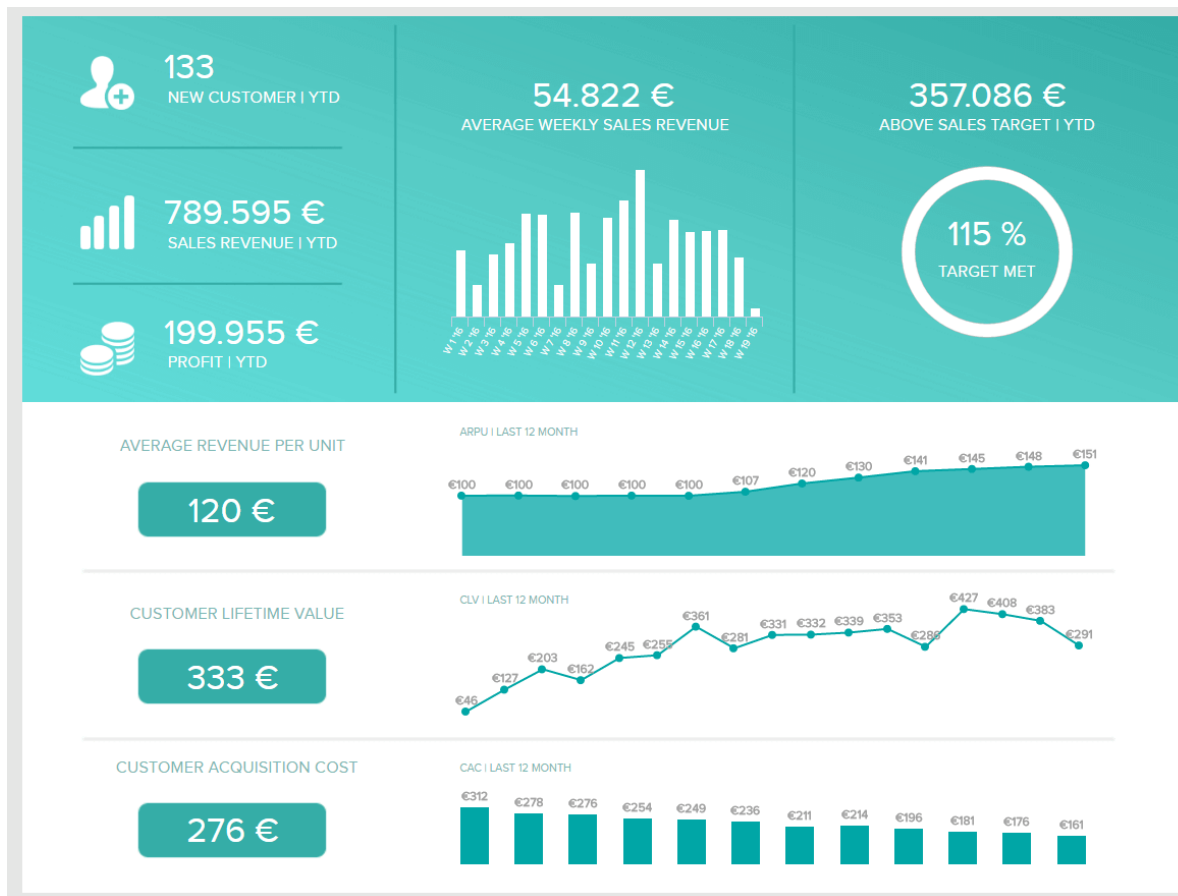
## | 6. *Visualisation des données*

La **data visualization** est la représentation graphique des données. Il s'agit de produire des images qui communiquent des informations abstraites pour les utilisateurs. Cette communication est réalisée par l'utilisation d'une cartographie.

Pour communiquer l'information de façon claire et efficace, la data visualisation utilise des graphiques statistiques, des diagrammes, des infographies et d'autres outils. Les données numériques peuvent être codées à l'aide de points, de lignes ou de barres pour communiquer visuellement un message quantitatif.

La data visualization permet de repérer des modèles, des tendances et des corrélations qui, autrement, passeraient inaperçus dans les rapports, tableaux ou feuilles de calcul traditionnels.

## 6. Visualisation des données





## | 7. *KPI - Tableau de bord*

- Le **suivi de gestion** est l'activité informationnelle continue de monitoring de chacune des activités opérationnelles qui forment les processus métiers implémentés dans le système d'information organisationnel ou l'OLTP, menant à des actions ou à des décisions d'ajustement régulier à chaque phase du processus.

- Un **KPI** ( Key performance Indicator - indicateur de performance) se définit comme étant une statistique ciblée et contextualisée selon une préoccupation de mesure, résultant de la collecte de données sur un élément lié au fonctionnement d'une organisation.

Le KPI constitue le maillon *INDISPENSABLE* du suivi de gestion puisqu'elle permet de mesurer le niveau d'accomplissement des activités qui forment le processus. Armé de la connaissance de cette mesure, les gestionnaires peuvent alors coordonner les activités du processus métier correctement.

- Les KPI sont regroupés de façon expressive dans ce que l'on appelle un **tableau de bord** et ensemble ceux-ci permettent de rapprocher le gestionnaire de l'information stratégique nécessaire au pilotage sain de son entreprise.

# La suite Microsoft BI

- Moteur de base de données : SQL Server Management Studio (**SSMS**)
- Outil d'intégration et de transformation des données (ETL) : SQL Server Integration Services (**SSIS**)
- Moteur multidimensionnel (cubes OLAP) et d'exploration des données : SQL Server Analysis Services (**SSAS**)
- Plateforme complète dédiée au Reporting : SQL Server Reporting Services (**SSRS**)
- Outil de reporting nouvelle génération : **Power BI**



*Questions ?*