

“Modeling ETL Activities as Graphs” (Vassiliadis, Simitsis, Skiadopoulos) y posicionamiento de la propuesta de “Primitivas de Transformación de Esquemas”

“Modeling ETL Activities...”

- Introduccion
- El “Architecture Graph”
- Explotación del Architecture graph
- ARKTOS
- Conclusiones

Adriana Marotta – jun/02

2

Introduccion

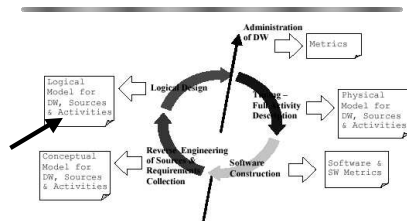


Fig. 1.2 The lifecycle of a Data Warehouse and its ETL processes

El trabajo se centra en el diseño lógico del escenario de ETL de un DW.

Adriana Marotta – jun/02

3

Introduccion

- Estado del arte:
 - Herramientas ETL
 - Ajax
 - Artículo de Workflow para ETL, de Mokrane
 - Otros trabajos de Vassiliadis sobre procesos de DW, metadata y calidad
- Pero...
 - Estos enfoques no han considerado en profundidad la *estructura interna de las actividades de ETL*.

Adriana Marotta – jun/02

4

Introduccion

- Contribuciones
 - Definición de un modelo lógico formal como abstracción lógica de procesos ETL
 - Incluye depositos de datos, actividades ETL y sus componentes (son todas entidades)
 - Reducción del modelo a un grafo, *Architecture Graph*
 - Las *entidades* anteriores son *nodos*, y las *relaciones* (4 tipos) entre ellas son las *aristas*.
 - Explotación del *Architecture Graph*
 - Transformaciones del grafo (zoom in/out)
 - Importance metrics: *dependence* y *responsibility*

Adriana Marotta – jun/02

5

El “Architecture Graph”

- Definiciones previas
- Escenario ETL
- Ejemplo
- Architecture Graph
- Ejemplo (cont.)

Adriana Marotta – jun/02

6

Definiciones previas

- **Nodos**
 - Tipos de datos
 - Tipos de funciones
 - Constantes
 - Atributos
 - Actividades
 - Conj. de registros
 - Funciones
- **Aristas (relaciones entre nodos)**
 - Part-of
 - Instance-of
 - Proveedor
 - Regulador
 - Proveedor derivado

Adriana Marotta – jun/02

7

Nodos

- **Data Type** - nombre, dominio
- **Attributes** - nombre, tipo de datos
- **RecordSet** - nombre, esquema, extension
- **Function Type** - nombre, tipos de datos de parametros, tipo de datos de resultado.
- **Function** - instancia de un tipo de funcion
- **Activities**
 - Terminología de [WfMC98] para procesos/programas
 - Abstracciones lógicas representando parte/modulos completos de código
 - Usan sentencias SQL para representar la semantica de las actividades

Adriana Marotta – jun/02

8

Aristas

- **Part-of** - Relaciona atributos y parámetros con actividades, conj. de registros o funciones
- **Instance-of** - Relacion entre data/function type y sus instancias
- **Provider** - Relaciones 1:N entre atributos (provider-consumer). Se define con nombre, mapping.
- **Regulator** - Relaciones entre los parametros de las actividades y los terminos que las instancian (populate).
- **Derived provider** - Caso particular de provider, cuando output attributes se derivan de la combinacion de input attributes y parametros.

Adriana Marotta – jun/02

9

Activities

- **Actividad, descrita formalmente por:**
 - *Name*
 - *Input Schema*
 - Recibe los datos del data provider
 - *Output Schema*
 - “placeholder” para las tuplas que pasan el chequeo de la activ.
 - *Rejections Schema*
 - “placeholder” para las tuplas que no pasan el chequeo de la activ.
 - *Parameter list* (esquema, atributo, funcion o cte.)
 - *Output Operational Semantics* (SQL)
 - *Rejection Operational Semantics* (por def. la neg. del ant.)

Adriana Marotta – jun/02

10

Escenario ETL

- **Consiste de**
 - Name
 - Activities
 - RecordSets (fuente)
 - Targets
 - tablas del DW que serán pobladas por las actividades
 - Provider Relationships
 - Relaciones “provider” entre actividades y recordsets del escenario
- **Intuitivamente**
 - Conj. de actividades presentadas en un grafo en una secuencia de ejecución.

Adriana Marotta – jun/02

11

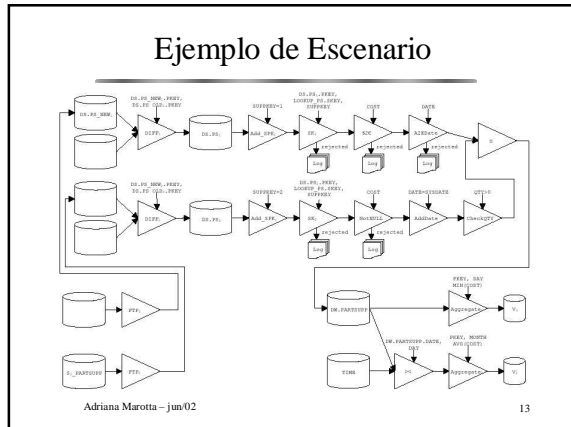
Ejemplo de Escenario

Source	Recordset Name	Recordset Schema
S ₁	S ₁ .PARTSUPP	PKKEY, DATE, QTY, COST
S ₂	S ₂ .PARTSUPP	PKKEY, QTY, COST
USA	DS.PS_NEW ₁	PKKEY, DATE, QTY, COST
	DS.PS_OLD ₁	PKKEY, DATE, QTY, COST
	DS.PS ₁	PKKEY, DATE, QTY, COST
	DS.PS_NEW ₂	PKKEY, QTY, COST
	DS.PS_OLD ₂	PKKEY, QTY, COST
	DS.PS ₂	PKKEY, QTY, COST
	DS.PS ₃	PKKEY, QTY, COST
DW	DW.PARTSUPP	PKKEY, SUPPKKEY, DATE, QTY, COST
	LOOKUP_PS	PKKEY, SOURCE, SKEY
	V1	PKKEY, DAY, MIN_COST
	V2	PKKEY, MONTH, AVG_COST
	TIME	DAY, MONTH, YEAR

Fig. 2.2 The schemata of the source databases and of the data warehouse

Adriana Marotta – jun/02

12



Architecture Graph








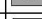




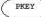
- El Architecture Graph y el escenario
 - El AG contiene todas las actividades y depositos de datos del escenario, junto con sus componentes
 - El AG captura el flujo de datos dentro del ambiente ETL
 - El AG da la informacion de
 - tipos de las entidades
 - regulacion de la ejecucion del escenario

Dominios y notacion

	Entity	Model-specific	Scenario-specific
Built-in	Data Types	D ¹	D
	Function Types	F ¹	F
	Constants	C ¹	C
User-provided	Attributes	Ω ¹	Ω
	Functions	Φ ¹	Φ
	Schemata	S ¹	S
	RecordSets	RS ¹	RS
	Activities	A ¹	A
	Provider Relationships	Pr ¹	Pr
	Part-Of Relationships	Po ¹	Po
	Instance-Of Relationships	Io ¹	Io
	Regulator Relationships	Rr ¹	Rr
	Derived Provider Relationships	Dr ¹	Dr

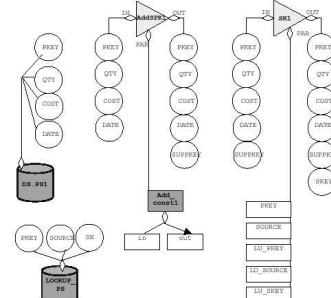
Fig. 3.1 Formal definition of domains and notation

Notacion grafica (AG)

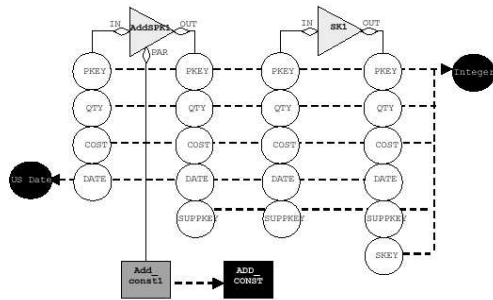
Data Types	Black ellipsis		RecordSets	Cylinders	
Function Types	Black squares		Functions	Gray squares	
Constants	Black circles		Parameters	White squares	
Attributes	Hollow ellipsoid nodes		Activities	Triangles	
Part-Of Relationships	Simple edges annotated with diamond*		Provider Relationships	Bold solid rows (from provider to consumer)	
Instance-Of Relationships	Dotted arrows (from instance towards the type)		Derived Provider Relationships	Bold dotted arrows (from provider to consumer)	
Regulator Relationships	Dotted edges		<p>* We annotate the part-of relationship among a function and its return type with a directed edge, to distinguish it from the rest of the consumers.</p>		

Ejemplo (cont.)

- Actividades: *Add_SPK_i* y *SK_i*
- Partes del grafo
 - Atributos y relaciones part-of
 - Data types y relaciones instance-of
 - Parametros y relaciones regulator
 - Relaciones provider
 - Relaciones derived provider



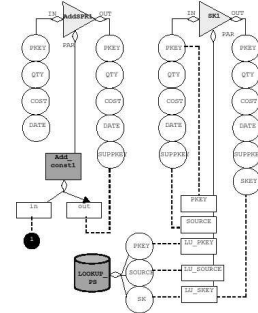
Data types y relaciones instance-of



Adriana Marotta - jun/02

19

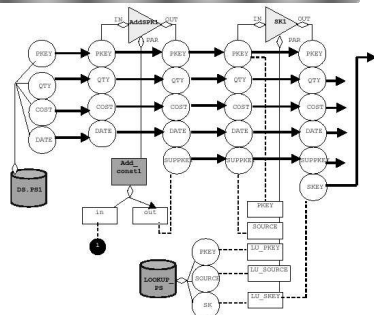
Parametros y relaciones regulator



Adriana Marotta - jun/02

20

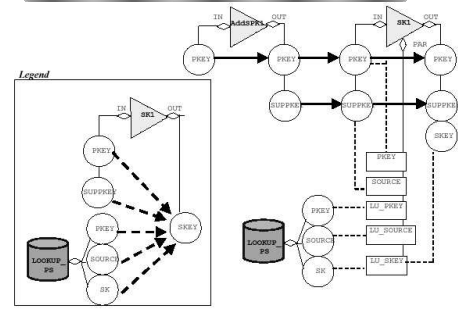
Relaciones provider



Adriana Marotta - jun/02

21

Relaciones derived provider



Adriana Marotta - jun/02

22

Explotación del AG

- Graph transformations
 - “zoom-out” para eliminar la sobrecarga de información (gran cantidad de atributos).
 - Subgrafo crítico que incluye solo las entidades necesarias para la población de los “target recordsets” del escenario.
- Importance metrics
 - Se asignan a los nodos para medir cuan crucial es su existencia para la ejecución del escenario
 - *local dependency, local responsibility, local degree*

Adriana Marotta - jun/02

23

Importance metrics

- Local dependency (in-degree)
 - Numero de nodos que deben ser activados para poblarlo
- Local responsibility (out-degree)
 - Cuantos nodos esperan por él para ser activados
- Local degree
 - La suma de los anteriores
- Transitividad
- Estas medidas no son aplicables solo a atributos. También a actividades y recordsets.

Adriana Marotta - jun/02

24

ARKTOS - Herramienta

- Objetivos
 - Facilidades graficas y declarativas para la definicion de tareas de limpieza y transformacion de DW
 - Medicion de calidad de datos (factores de calidad)
 - Ejecucion optimizada de secuencias complejas de tareas
- Metamodelo
 - El presentado anteriormente + *Error Type, Policy, Quality Factors*
- *Transformation and cleaning primitives*

Adriana Marotta – jun/02

25

Primitivas

- Operaciones primitivas para soportar el proceso ETL
- Personalizadas por el usuario
 - Input, output, data stores, policy, quality factors

Primitive Operation	SQL statement	SEMANTIC'S clause shorthand
UNIQUENESS VIOLATION	SELECT * FROM <table> GROUP BY <attribute> HAVING COUNT(*) > 1	<table>.<attribute>
NULL REFERENCE	SELECT * FROM <table> WHERE <attribute> IS NULL	<table>.<attribute>
DOMAIN VIOLATION	SELECT * FROM <table> WHERE <attribute> NOT IN <domain specification>	<table>.<attribute>
PRIMARY KEY VIOLATION	SELECT * FROM <table> GROUP BY (<attribute_1>,...,<attribute_n>) HAVING COUNT(*) > 1	<table>.<attribute_1>,...,<attribute_n>
REFERENCE VIOLATION	SELECT * FROM <table> WHERE <attribute> NOT IN (SELECT <target_attribute> FROM <target_table>)	<table>.<attribute> NOT IN <target_table>.<target_attribute>
FORMAT VIOLATION	SELECT APPLY(<regexp>,<attribute>) FROM <table> WHERE APPLY(<regexp>,<attribute>)	TARGET APPLY(<regexp>,<attribute>) SOURCE APPLY(<regexp>,<attribute>)
NULL	Ad-hoc SQL query	Arbitrary SQL query

— where <reg_exp> is PERL regular expression acting as a formatting function

Adriana Marotta – jun/02

26

Definicion de actividades

- Graficamente
 - Paleta con las actividades provistas por ARKTOS
- Declarativamente
 - 2 lenguajes
 - XADL
 - XML-based Activity Definition Language
 - SADL
 - Simple Activity Definition Language

Adriana Marotta – jun/02

27

Conclusiones – [VASS02]

- Diseño logico de escenario ETL
- Modelo logico formal y “Architecture Graph”
- “Importance metrics”
- Herramienta ARKTOS para modelar y ejecutar escenarios
- Primitivas para las tareas comunes
- Resultados son parte de un proyecto mas grande
 - Referencian: “Conceptual modeling for ETL Processes”

Adriana Marotta – jun/02

28

Trabajo futuro – [VASS02]

- Optimizacion
- Proveer primitivas de transformacion mas ricas
- “It would be nice if the research community could provide formal design guidelines (in the sense of normal forms or extra integrity constraints) for the engineering of ETL processes.”

Adriana Marotta – jun/02

29

Frases célebres – [VASS02]

- Sobre VM
 - “... viewing the DW as a set of layered, materIALIZED views is a very simplistic view.”
 - “... insufficient to describe the structure and contents of a DW.”
- Sobre evolucion
 - “... the interesting problem is how to design the scenario in order to achieve effectiveness, efficiency and tolerance of the impacts of evolution.”
 - “Dependence and responsibility are crucial measures for the engineering of the evolution of the ETL env.”

Adriana Marotta – jun/02

30

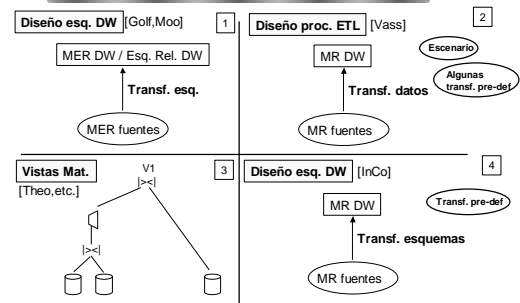
Posicionamiento de las Primitivas de Transformacion de Esquemas

- Algunos enfoques... y el nuestro
- Como nos podemos posicionar?
- Ejemplo
- Conclusiones ?

Adriana Marotta – jun/02

31

Algunos enfoques... y el nuestro



Adriana Marotta – jun/02

32

Cómo nos podemos posicionar?

- Diagnóstico
 - Las propuestas 1 y 2 estan totalmente desconectadas.
 - La propuesta de diseño no se preocupa por los procesos ETL y la propuesta para ETL no menciona su relacion con el diseño del esquema de DW, por ejemplo, cómo deduce las actividades a realizar?
 - Las propuestas tipo 1 no resuelven la construccion de estructuras “complejas” de DW (historicos, dimension versioning, datos calculados, etc.), que a la vez son comunes.

Adriana Marotta – jun/02

33

Cómo nos podemos posicionar?

- “Filling the gap between DW schema design and the corresponding ETL processes design”
- Un mecanismo para diseño logico de DW que ofrece 2 ventajas:
 - Facilidades para diseñar estructuras complejas de DW
 - Permite la deducccion semi-automatica de un “Architecture Graph” [Vass].

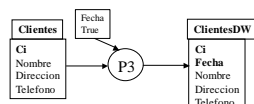
Adriana Marotta – jun/02

34

Ejemplo

- Una de las estrategias de diseño para Dimension Versioning era aplicar *Temporalization* haciendo que el nuevo atributo forme parte de la clave.

Diseño del esquema:

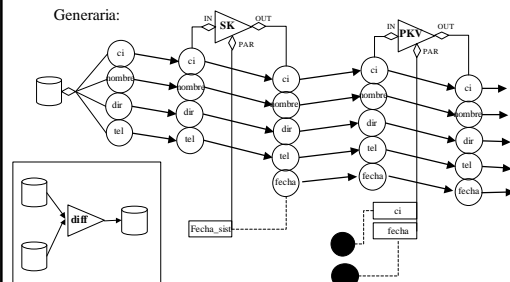


Adriana Marotta – jun/02

35

Ejemplo cont.

Generaria:



Adriana Marotta – jun/02

36

Conclusiones ?
