

## Data Integration Exercise

Technologies for Information Systems  
November, 7<sup>th</sup> 2011



Politecnico di Milano  
Dipartimento di Elettronica e Informazione



### Problem setting

MediaWiki and TikiWiki are two wikilike content management systems with a DB backend.

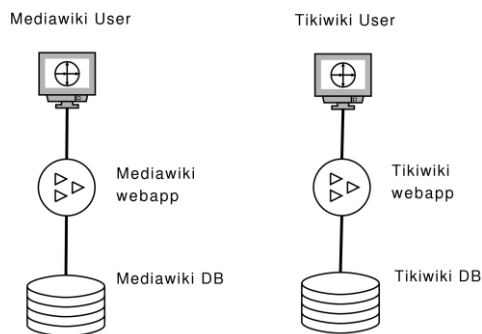
At the University of Wasomao some of the professors use MediaWiki while some others use TikiWiki.

The university wants to allow the professor to use one of the two systems to access the pages of the other one, while keeping unaltered the original applications.

TSI 10-11



### Current Situation



TSI 10-11



### Original DB backends

#### MEDIAWIKI-DB

**PAGE** (pid, title, namespace)

**REVISION** (rid, page, text-id)

**TEXT**(tid, plain-text)

**USER** (uid, nickname, password);

**PAGELINKS** (from, title, namespace)

*PAGE* contains the metadata about a page

*REVISION* contains different versions of the same page

*TEXT* contains individual text of each revision.

*USER* contains the information about CMS users

*PAGELINKS* stores the links between pages as an id of the source page and the title and namespace of the target page.

TSI 10-11



### Original DB backends (2)

#### TIKIWIKI-DB

**PAGES** (page-id, version, page-title, user-id, text)

**COMMENTS** (page-id, user-id, comment)

**USER** (user-id, real-name, username, email, pwd, path-to-picture)

**LINKS** (from-page, to-page)

*PAGES* contains the information about pages (their versions and the actual texts)

*COMMENTS* stores user comments to pages, *USER* contains information about users

*LINKS* represents the links between pages as id pairs of the source page and of the target page.

TSI 10-11



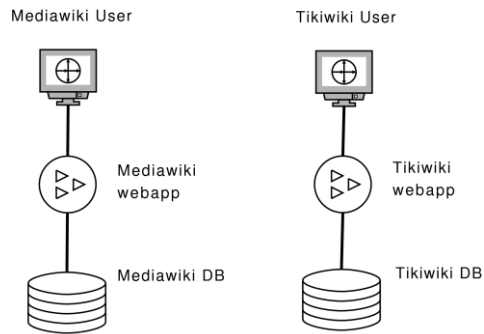
### Exercise

- 1) Propose a data-integration solution based on the requirements above
- 2) Provide, for each input data source, the reverse engineering from the logical to conceptual schema (ER graph)
- 3) Design the mappings between the two sources
- 4) Consider the following query:  
*Return all the pages pointing to a page titled "data integration"*
  - Show how this query will be rewritten when posed on both systems

TSI 10-11



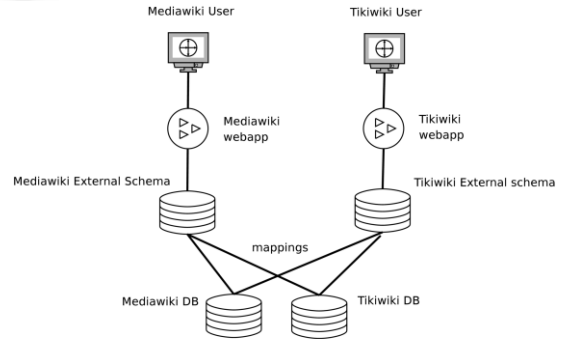
## Current Situation



TSI 10-11



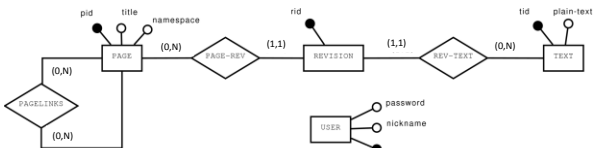
## Proposed solution



TSI 10-11



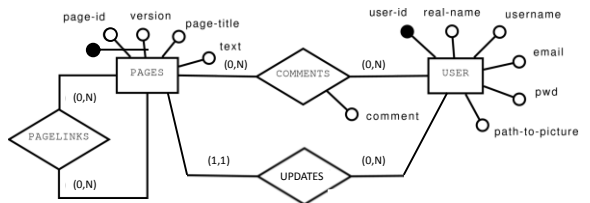
## Mediawiki reverse engineering



TSI 10-11



## Tikiwiki reverse engineering



TSI 10-11



## Mapping design

To ease the query processing we design two set of mappings:

- TIKITOMEDIA which considers the MEDIAWIKI schema as the target schema and the TIKIWIKI schema as the source schema.
- MEDIATOTIKI which does the opposite.

TSI 10-11



## TIKITOMEDIA mappings

```
CREATE VIEW PAGE (pid, title, namespace) AS
SELECT      *
FROM        MEDIA.PAGE

UNION

SELECT      DISTINCT TP.page-id, TP.page-title, 'main'
FROM        TIKI.PAGES as TP

CREATE VIEW TEXT (tid, plain-text) AS
SELECT      *
FROM        MEDIA.TEXT

UNION

SELECT      ENCRYPT_TDES(TP.text, 'abc'), TP.text
FROM        TIKI.PAGE as TP
```

TSI 10-11



## TIKITOMEDIA mappings (2)

```
CREATE VIEW REVISION (rid, page, text-id) AS

SELECT      *
FROM        MEDIA.REVISION

UNION

SELECT      CONCAT(TP.page-id, '$', TP.version),
            TP.page-id, ENCRYPT_TDES(TP.text, 'abc')
FROM        FROM TIKI.PAGES AS TP

CREATE VIEW USER (uid, nickname, password) AS

SELECT      *
FROM        MEDIA.USER

UNION

SELECT      TU.user-id, TU.username, TU.pwd
FROM        TIKI.USER
```

TSI 10-11



## TIKITOMEDIA mappings (3)

```
CREATE VIEW PAGELINKS (from, title, namespace) AS

SELECT      *
FROM        MEDIA.PAGELINKS

UNION

SELECT      TL.from-page, TP.page-title, 'main'
FROM        TIKI.LINKS AS TL, TIKI.PAGES AS TP
WHERE       TP.page-id = TL.to-page
```

TSI 10-11



## MEDIATOTIKI mappings

```
CREATE VIEW PAGES (page-id, version, page-title, user_id, text) AS

SELECT      *
FROM        TIKI.PAGES

UNION

SELECT      MP.pid, MR.rid, MP.title, null, MT.plain-text
FROM        MEDIA.PAGE as MP, MEDIA.REVISION as MR,
            MEDIA.TEXT as MT
WHERE       MP.pid = MR.page and MR.text-id = MT.tid

CREATE VIEW USER (user-id, real-name, username, email, pwd, path-to-
picture) AS

SELECT      *
FROM        TIKI.USER

UNION

SELECT      MU.uid, null, MU.nickname, null, MU.password, null
FROM        MEDIA.USER as MU
```

TSI 10-11



## MEDIATOTIKI mappings (2)

```
CREATE VIEW LINKS (from-page, to-page) AS

SELECT      *
FROM        TIKI.LINKS

UNION

SELECT      MPL.from, MP.p-id
FROM        MEDIA.PAGELINKS as MPL, MEDIA.PAGE as MP
WHERE       MPL.title = TP.title
            AND MPL.namespace = MP.namespace

CREATE VIEW COMMENTS (page-id, user-id, comment) AS

SELECT      *
FROM        TIKI.COMMENTS
```

TSI 10-11



## Query

Return all the pages pointing to a page titled "data integration"

TSI 10-11



## Query Rewriting - Tikiwiki external schema - GAV

*Tikiwiki*

```
Q(x) ← LINKS(x,y) ∧ PAGES(y,_,w,_) ∧ w='data integration'
where
v1: LINKS(x,y) ← PAGELINKS(x,z,q) ∧ PAGE(y,z,q)
v2: PAGES (x,y,z,j,v) ← PAGE(x,z,_) ∧ REVISION(y,x,s) ∧ TEXT (s,v)
```

TSI 10-11



## Query Rewriting - Tikiwiki - GAV (2)

Tikiwiki

$Q(x) \leftarrow \text{LINKS}(x,y) \wedge \text{PAGES}(y,_,w,_,_) \wedge w = \text{'data integration'}$

where

$v1: \text{LINKS}(x,y) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q)$

$v2: \text{PAGES}(x,y,z,j,v) \leftarrow \text{PAGE}(x,z,_) \wedge \text{REVISION}(y,x,s) \wedge \text{TEXT}(s,v)$

rewritten query (simple unfolding – GAV mappings)

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{PAGE}(y,w,_) \wedge \text{REVISION}(_,y,s) \wedge \text{TEXT}(s,_) \wedge w = \text{'data integration'}$

TSI 10-11



## Query Rewriting - Tikiwiki - GAV (3)

Tikiwiki

$Q(x) \leftarrow \text{LINKS}(x,y) \wedge \text{PAGES}(y,_,w,_,_) \wedge w = \text{'data integration'}$

where

$v1: \text{LINKS}(x,y) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q)$

$v2: \text{PAGES}(x,y,z,j,v) \leftarrow \text{PAGE}(x,z,_) \wedge \text{REVISION}(y,x,s) \wedge \text{TEXT}(s,v)$

rewritten query (simple unfolding – GAV mappings)

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{PAGE}(y,w,_) \wedge \text{REVISION}(_,y,s) \wedge \text{TEXT}(s,_) \wedge w = \text{'data integration'}$

Unification possible because y is the primary key ( $w \rightarrow z$ )

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{REVISION}(_,y,s) \wedge \text{TEXT}(s,_) \wedge z = \text{'data integration'}$

TSI 10-11



## Query Rewriting - Tikiwiki - GAV (4)

Tikiwiki

$Q(x) \leftarrow \text{LINKS}(x,y) \wedge \text{PAGES}(y,_,w,_,_) \wedge w = \text{'data integration'}$

where

$v1: \text{LINKS}(x,y) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q)$

$v2: \text{PAGES}(x,y,z,j,v) \leftarrow \text{PAGE}(x,z,_) \wedge \text{REVISION}(y,x,s) \wedge \text{TEXT}(s,v)$

rewritten query (simple unfolding – GAV mappings)

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{PAGE}(y,w,_) \wedge \text{REVISION}(_,y,s) \wedge \text{TEXT}(s,_) \wedge w = \text{'data integration'}$

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{REVISION}(_,y,s) \wedge \text{TEXT}(s,_) \wedge z = \text{'data integration'}$

We know from the ER schema that there is a foreign key constraint from REVISION to TEXT. Therefore the TEXT(s,\_) atom can be deleted, since it does not impose any extra condition (the variables in the atom are not used in any other atom apart the two atoms involved in the foreign key constraint), if we assume the data to be consistent with the constraints.

TSI 10-11



## Query Rewriting - Tikiwiki - GAV (5)

Tikiwiki

$Q(x) \leftarrow \text{LINKS}(x,y) \wedge \text{PAGES}(y,_,w,_,_) \wedge w = \text{'data integration'}$

where

$v1: \text{LINKS}(x,y) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q)$

$v2: \text{PAGES}(x,y,z,j,v) \leftarrow \text{PAGE}(x,z,_) \wedge \text{REVISION}(y,x,s) \wedge \text{TEXT}(s,v)$

rewritten query (simple unfolding – GAV mappings)

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{PAGE}(y,w,_) \wedge \text{REVISION}(_,y,s) \wedge \text{TEXT}(s,_) \wedge w = \text{'data integration'}$

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{REVISION}(_,y,s) \wedge \text{TEXT}(s,_) \wedge z = \text{'data integration'}$

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{REVISION}(_,y,s) \wedge z = \text{'data integration'}$

TSI 10-11



## Query Rewriting - Tikiwiki - GAV (6)

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{REVISION}(_,y,s) \wedge z = \text{'data integration'}$

Similarly, we remove the REVISION(.,y,s) atom since there must be a revision related to each page (as we can see from the ER schema).

TSI 10-11



## Query Rewriting - Tikiwiki - GAV (7)

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{REVISION}(_,y,s) \wedge z = \text{'data integration'}$

Similarly, we remove the REVISION(.,y,s) atom since there must be a revision related to each page (as we can see from the ER schema).

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge z = \text{'data integration'}$

TSI 10-11



## Query Rewriting - Tikiwiki - GAV (8)

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{REVISION}(\_,y,s) \wedge z = \text{'data integration'}$

Similarly, we remove the  $\text{REVISION}(\_,y,s)$  atom since there must be a revision related to each page (as we can see from the ER schema).

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge z = \text{'data integration'}$

Although this is not in the ER schema, it might be reasonable to assume that the following inclusion dependency holds:

$\text{PAGELINKS}[\text{title}, \text{namespaces}] \subseteq \text{PAGE}[\text{title}, \text{namespaces}]$

TS1 10-11



## Query Rewriting - Tikiwiki - GAV (9)

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge \text{REVISION}(\_,y,s) \wedge z = \text{'data integration'}$

Similarly, we remove the  $\text{REVISION}(\_,y,s)$  atom since there must be a revision related to each page (as we can see from the ER schema).

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge \text{PAGE}(y,z,q) \wedge z = \text{'data integration'}$

Although this is not in the ER schema, it might be reasonable to assume that the following inclusion dependency holds:

$\text{PAGELINKS}[\text{title}, \text{namespaces}] \subseteq \text{PAGE}[\text{title}, \text{namespaces}]$

If so, for the same reasons as before, we can simplify the query as follows:

$Q(x) \leftarrow \text{PAGELINKS}(x,z,q) \wedge z = \text{'data integration'}$

TS1 10-11



## Query Rewriting - Mediawiki external schema-GAV

*Mediawiki*

$Q(x) \leftarrow \text{PAGELINKS}(x,y,\_) \wedge y = \text{'data integration'}$

where

$v1: \text{PAGELINKS}(x,y,\text{'main'}) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_)$

TS1 10-11



## Query Rewriting - Mediawiki-GAV(1)

*Mediawiki*

$Q(x) \leftarrow \text{PAGELINKS}(x,y,\_) \wedge y = \text{'data integration'}$

where

$v1: \text{PAGELINKS}(x,y,\text{'main'}) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_)$

We rewrite the view in such a way that the head only contains distinct variables:

$v1: \text{PAGELINKS}(x,y,z) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_) \wedge z = \text{'main'}$

TS1 10-11



## Query Rewriting - Mediawiki-GAV(2)

*Mediawiki*

$Q(x) \leftarrow \text{PAGELINKS}(x,y,\_) \wedge y = \text{'data integration'}$

where

$v1: \text{PAGELINKS}(x,y,\text{'main'}) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_)$

We rewrite the view in such a way that the head only contains distinct variables:

$v1: \text{PAGELINKS}(x,y,z) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_) \wedge z = \text{'main'}$

rewritten query (simple unfolding – GAV mappings)

$Q(x) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_) \wedge y = \text{'data integration'} \wedge \_ = \text{'main'}$

TS1 10-11



## Query Rewriting - Mediawiki-GAV(2)

*Mediawiki*

$Q(x) \leftarrow \text{PAGELINKS}(x,y,\_) \wedge y = \text{'data integration'}$

where

$v1: \text{PAGELINKS}(x,y,\text{'main'}) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_)$

We rewrite the view in such a way that the head only contains distinct variables:

$v1: \text{PAGELINKS}(x,y,z) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_) \wedge z = \text{'main'}$

rewritten query (simple unfolding – GAV mappings)

$Q(x) \leftarrow \text{LINKS}(x,w) \wedge \text{PAGES}(w,\_,y,\_) \wedge y = \text{'data integration'} \wedge \_ = \text{'main'}$

In the last atom there is a variable ( $\_$ ) that is not used elsewhere, so we can safely remove it:

TS1 10-11



## Query Rewriting - Mediawiki-GAV(2)

### *Mediawiki*

$Q(x) \leftarrow \text{PAGELINKS}(x, y, \_) \wedge y = \text{'data integration'}$

where

$v1: \text{PAGELINKS}(x, y, \text{'main'}) \leftarrow \text{LINKS}(x, w) \wedge \text{PAGES}(w, \_, y, \_)$

We rewrite the view in such a way that the head only contains distinct variables:

$v1: \text{PAGELINKS}(x, y, z) \leftarrow \text{LINKS}(x, w) \wedge \text{PAGES}(w, \_, y, \_) \wedge z = \text{'main'}$

rewritten query (simple unfolding – GAV mappings)

$Q(x) \leftarrow \text{LINKS}(x, w) \wedge \text{PAGES}(w, \_, y, \_) \wedge y = \text{'data integration'} \wedge \_ = \text{'main'}$

In the last atom there is a variable ( $\_$ ) that is not used elsewhere, so we can safely remove it:

$Q(x) \leftarrow \text{LINKS}(x, w) \wedge \text{PAGES}(w, \_, y, \_) \wedge y = \text{'data integration'}$