# Lecture 2
# Entity-Relationship Model

Eugene Wu

# Steps for a New Application

Requirements
   what are you going to build?
Conceptual Database Design
   pen-and-pencil description
Logical Design
   formal database schema
Schema Refinement:
   fix potential problems, normalization
Physical Database Design
   use sample of queries to optimize for speed/storage
App/Security Design
   prevent security problems

# Steps for a New Application

Requirements
  what are you going to build?

Conceptual Database Design                    ER Modeling
  pen-and-pencil description

Logical Design
  formal database schema

Schema Refinement:
  fix potential problems, normalization

Physical Database Design
  use sample of queries to optimize for speed/storage

App/Security Design
  prevent security problems

# Database Apps Are Complicated

Typical Fortune 100 Company

    ~10k different information (data) systems

    90% relational databases (DBMSes)
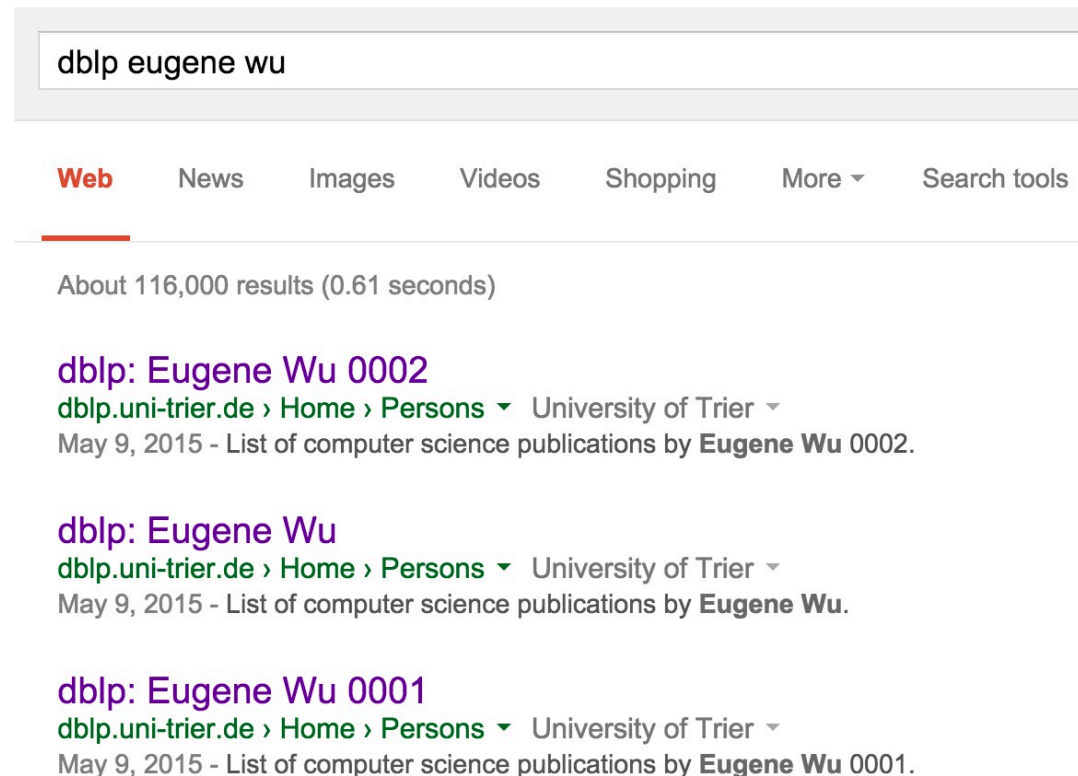
    Typical database has >100 tables

    Typical table has 50 – 200 attributes

# Inconsistencies/Constraint Violations

*Huge* amount of effort to avoid inconsistencies

Can data model help us avoid automatically?

DBLP is *the* site for computer science publications

# Inconsistencies/Constraint Violations

**2014**

■ [j8]   Eugene Wu, Leilani Battle, Samuel R. Madden:
**The Case for Data Visualization Management Systems.** PVLDB 7(10): 903-906 (2014)

■ [j7]   Alekh Jindal, Praynaa Rawlani, Eugene Wu, Samuel Madden, Amol Deshpande, Mike Stonebraker:
**VERTEXICA: Your Relational Friend for Graph Analytics!** PVLDB 7(13): 1669-1672 (2014)

$$\neq$$

**1994**

■ [c2]   James Hwang, Eugene Wu, Alan Bell, Andy Cordell, LeBarian Stokes, Scott Hankins:
**Design of a SPDM-Like Robotic Manipulator System for Space Station on Orbit Replaceable Unit Ground Testing - An Overview of the System Architecture.** ICRA 1994: 1286-1291

■ [c1]   Eugene Wu, James Hwang, Scott Hankins:
**Design of the Control System for a Robotic Manipulator for Space Station On-Orbit Replaceable Unit Ground Testing.** ICRA 1994: 1415-1420

Not Teaching 4111!!!

# Inconsistencies/Constraint Violations

## Check in application code!

# It is Hard to Design Applications

Server Code

DBMS

# It is Hard to Design Applications

Add new users

Server Code

DBMS

# It is Hard to Design Applications

# It is Hard to Design Applications

# ER Diagrams

What is it?

- A way to sketch the core information that your database will eventually store.

- Visually encodes important constraint information

Who cares?

- Good for "white boarding" together

- Good way to share the "gist" of your DB's structure

```
test=# \d election
                          Table "public.election"
      Column       |  Type   | Collation | Nullable |            Default
-------------------+---------+-----------+----------+------------------------------
 year              | integer |           |          |
 state             | text    |           |          |
 state_po          | text    |           |          |
 state_fips        | integer |           |          |
 state_cen         | integer |           |          |
 state_ic          | integer |           |          |
 office            | text    |           |          |
 candidate         | text    |           |          |
 party_detailed    | text    |           |          |
 writein           | text    |           |          |
 candidatevotes    | integer |           |          |
 totalvotes        | integer |           |          |
 version           | integer |           |          |
 notes             | text    |           |          |
 party_simplified  | text    |           |          |
 id                | integer |           | not null | nextval('election_id_seq'::
Indexes:
    "election_id_key" UNIQUE CONSTRAINT, btree (id)

test=# \d food
                        Table "public.food"
      Column      |  Type   | Collation | Nullable | Default
------------------+---------+-----------+----------+---------
 camis            | integer |           |          |
 dba              | text    |           |          |
 boro             | text    |           |          |
 building         | integer |           |          |
 street           | text    |           |          |
 zipcode          | integer |           |          |
 phone            | bigint  |           |          |
 inspection_date  | date    |           |          |
 action           | text    |           |          |
 score            | integer |           |          |
 grade            | text    |           |          |
 inspection_type  | text    |           |          |
 census_tract     | integer |           |          |
 year             | integer |           |          |
 month            | integer |           |          |
 day              | integer |           |          |
```

Database schema of MediaWiki 1.24.1 (December 2014)
Refer to https://www.mediawiki.org/wiki/DB for more details.

## Users

**user_properties**
- up_user INT
- up_property VARBINARY(255)
- up_value BLOB
- Indexes

**user_groups**
- ug_user INT
- ug_group VARBINARY(255)
- Indexes

**user_former_groups**
- ufg_user INT
- ufg_group VARBINARY(255)
- Indexes

**user_newtalk**
- user_id INT
- user_ip VARBINARY(40)
- user_last_timestamp VARBINARY(14)
- Indexes

**user**
- user_id INT
- user_name VARCHAR(255)
- user_real_name VARCHAR(255)
- user_password TINYBLOB
- user_newpassword TINYBLOB
- user_newpass_time BINARY(14)
- user_email TINYTEXT
- user_touched BINARY(14)
- user_token BINARY(32)
- user_email_authenticated BINARY(14)
- user_email_token BINARY(32)
- user_email_token_expires BINARY(14)
- user_registration BINARY(14)
- user_editcount INT
- user_password_expires VARBINARY(14)
- Indexes

**ipblocks**
- ipb_id INT
- ipb_address TINYBLOB
- ipb_user INT
- ipb_by INT
- ipb_by_text VARCHAR(255)
- ipb_reason TINYBLOB
- ipb_timestamp BINARY(14)
- ipb_auto
- ipb_anon_only
- ipb_create_account
- ipb_enable_autoblock
- ipb_expiry VARBINARY(14)
- ipb_range_start TINYBLOB
- ipb_range_end TINYBLOB
- ipb_deleted
- ipb_block_email
- ipb_allow_usertalk
- ipb_parent_block_id INT
- Indexes

## Logging

**logging**
- log_id INT
- log_type VARBINARY(32)
- log_action VARBINARY(32)
- log_timestamp BINARY(14)
- log_user INT
- log_user_text VARCHAR(255)
- log_namespace INT
- log_title VARCHAR(255)
- log_page INT
- log_comment VARCHAR(255)
- log_params BLOB
- log_deleted TINYINT
- Indexes

**log_search**
- ls_field VARBINARY(32)
- ls_value VARCHAR(255)
- ls_log_id INT
- Indexes

## Tags

**change_tag**
- ct_rc_id INT
- ct_log_id INT
- ct_rev_id INT
- ct_tag VARCHAR(255)
- ct_params BLOB
- Indexes

**valid_tag**
- vt_tag VARCHAR(255)
- Indexes

**tag_summary**
- ts_rc_id INT
- ts_log_id INT
- ts_rev_id INT
- ts_tags BLOB
- Indexes

## Recent changes

**recentchanges**
- rc_id INT
- rc_timestamp VARBINARY(14)
- rc_user INT
- rc_user_text VARCHAR(255)
- rc_namespace INT
- rc_title VARCHAR(255)
- rc_comment VARCHAR(255)
- rc_minor TINYINT
- rc_bot TINYINT
- rc_new TINYINT
- rc_cur_id INT
- rc_this_oldid INT
- rc_last_oldid INT
- rc_type TINYINT
- rc_source VARCHAR(16)
- rc_patrolled TINYINT
- rc_ip VARBINARY(40)
- rc_old_len INT
- rc_new_len INT
- rc_deleted TINYINT
- rc_logid INT
- rc_log_type VARBINARY(255)
- rc_log_action VARBINARY(255)
- rc_params BLOB
- Indexes

**watchlist**
- wl_user INT
- wl_namespace INT
- wl_title VARCHAR(255)
- wl_notificationtimestamp VARBINARY(14)
- Indexes

## Pages

**archive**
- ar_id INT
- ar_namespace INT
- ar_title VARCHAR(255)
- ar_text MEDIUMBLOB
- ar_comment TINYBLOB
- ar_user INT
- ar_user_text VARCHAR(255)
- ar_timestamp BINARY(14)
- ar_minor_edit TINYINT
- ar_flags TINYBLOB
- ar_rev_id INT
- ar_text_id INT
- ar_deleted TINYINT
- ar_len INT
- ar_page_id INT
- ar_parent_id INT
- ar_sha1 VARBINARY(32)
- ar_content_model VARBINARY
- ar_content_format VARBINARY

**revision**
- rev_id INT
- rev_page INT
- rev_text_id INT
- rev_comment TINYBLOB
- rev_user INT
- rev_user_text VARCHAR(255)
- rev_timestamp BINARY(14)
- rev_minor_edit TINYINT
- rev_deleted TINYINT
- rev_len INT
- rev_parent_id INT
- rev_sha1 VARBINARY(32)
- rev_content_model VARBINARY(32)
- rev_content_format VARBINARY(64)
- Indexes

**page**
- page_id INT
- page_namespace INT
- page_title VARCHAR(255)
- page_restrictions TINYBLOB
- page_counter BIGINT
- page_is_redirect TINYINT
- page_is_new TINYINT
- page_random DOUBLE
- page_touched BINARY(14)
- page_links_updated VARBINARY(14)
- page_latest INT
- page_len INT
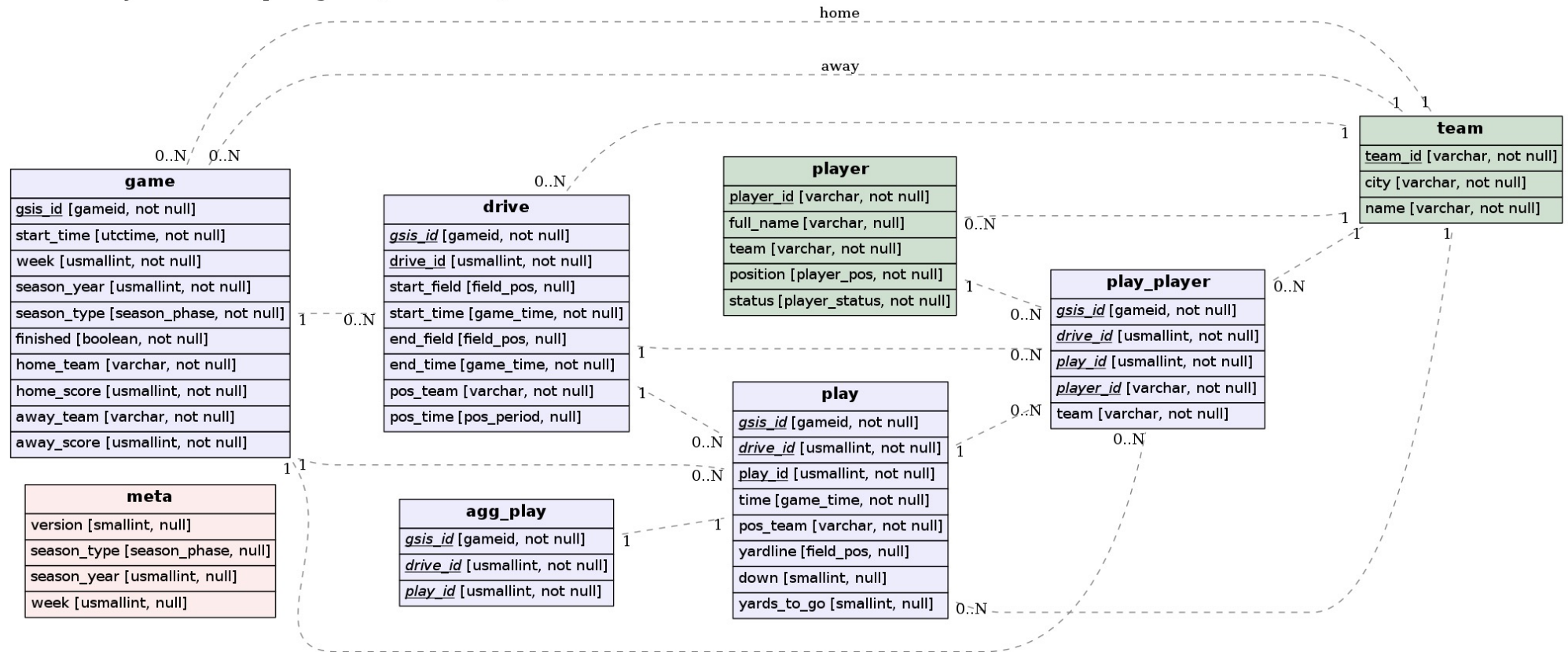- page_content_model VARBINARY(32)
- page_lang VARBINARY(35)
- Indexes

**page_props**
- pp_page INT
- pp_propname VARBINARY(60)
- pp_value BLOB
- pp_sortkey FLOAT
- Indexes

**page_restrictions**
- pr_id INT
- pr_page INT
- pr_type VARBINARY(60)
- pr_level VARBINARY(60)
- pr_cascade TINYINT
- pr_user INT
- pr_expiry VARBINARY(14)
- Indexes

**text**
- old_id INT
- old_text MEDIUMBLOB
- old_flags TINYBLOB
- Indexes

**redirect**
- rd_from INT
- rd_namespace INT
- rd_title VARCHAR(255)
- rd_interwiki VARCHAR(32)
- rd_fragment VARCHAR(255)
- Indexes

**category**
- cat_id INT
- cat_title VARCHAR(255)
- cat_pages INT
- cat_subcats INT
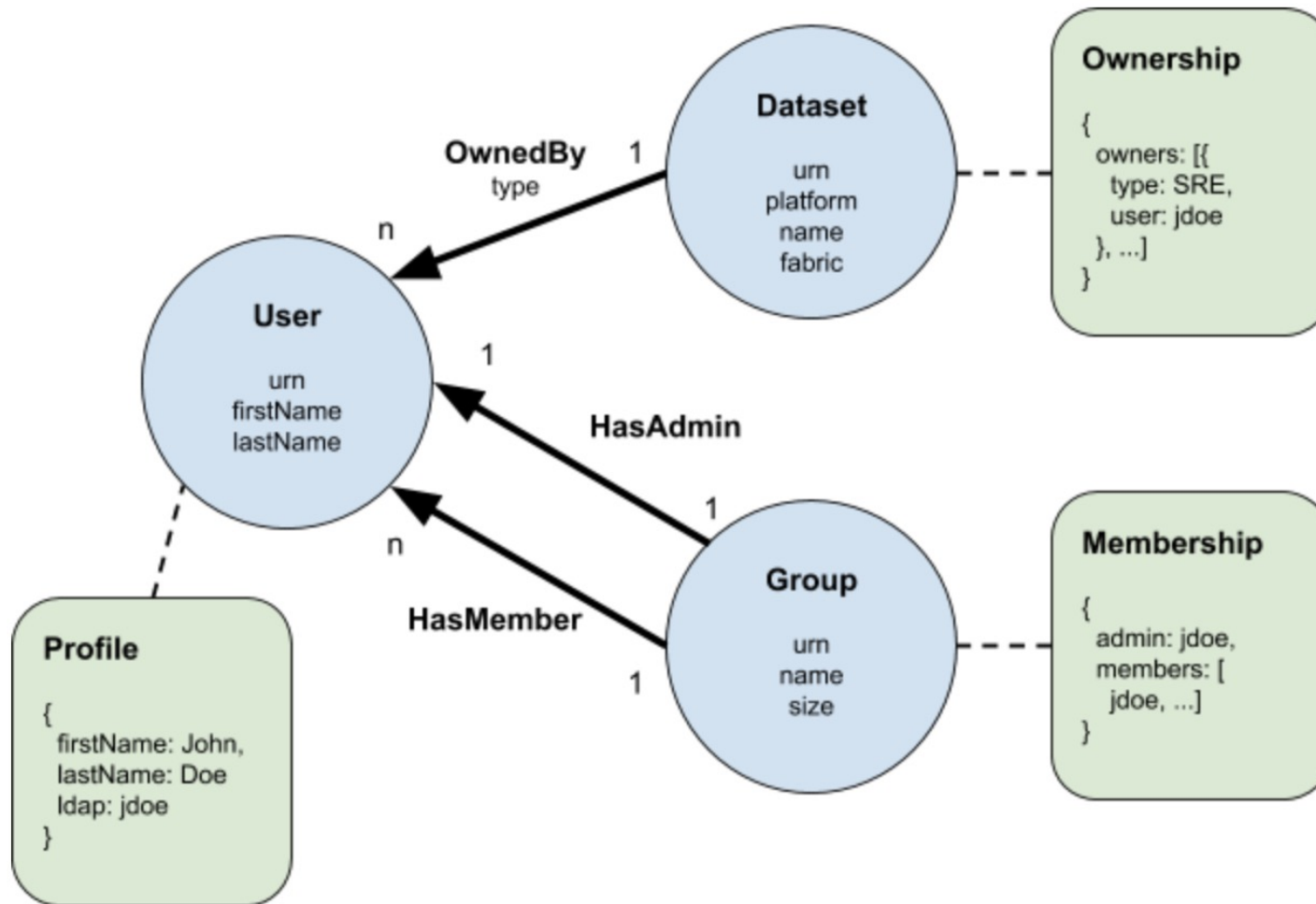- cat_files INT
- Indexes

**protected_titles**
- pt_namespace INT
- pt_title VARCHAR(255)
- pt_user INT
- pt_reason TINYBLOB
- pt_timestamp BINARY(14)
- pt_expiry VARBINARY(14)
- pt_create_perm VARBINARY(60)
- Indexes

## Link tables

**pagelinks**
- pl_from INT
- pl_from_namespace INT
- pl_namespace INT
- pl_title VARCHAR(255)
- Indexes

**imagelinks**
- il_from INT
- il_from_namespace INT
- il_to VARCHAR(255)
- Indexes

**iwlinks**
- iwl_from INT
- iwl_prefix VARBINARY(20)
- iwl_title VARCHAR(255)
- Indexes

**templatelinks**
- tl_from INT
- tl_from_namespace INT
- tl_namespace INT
- tl_title VARCHAR(255)
- Indexes

**externallinks**
- el_id INT
- el_from INT
- el_to BLOB
- el_index BLOB
- Indexes

**categorylinks**
- cl_from INT
- cl_to VARCHAR(255)
- cl_sortkey VARBINARY(230)
- cl_sortkey_prefix VARCHAR(255)
- cl_timestamp TIMESTAMP
- cl_collation VARBINARY(32)
- cl_type ENUM(...)
- Indexes

**langlinks**
- ll_from INT
- ll_lang VARBINARY(20)
- ll_title VARCHAR(255)
- Indexes

## Multimedia

**image**
- img_name VARCHAR(255)
- img_size INT
- img_width INT
- img_height INT
- img_metadata MEDIUMBLOB
- img_bits INT
- img_media_type ENUM(...)
- img_major_mime ENUM(...)
- img_minor_mime VARBINARY(100)
- img_description TINYBLOB
- img_user INT
- img_user_text VARCHAR(255)
- img_timestamp VARBINARY(14)
- img_sha1 VARBINARY(32)
- Indexes

**filearchive**
- fa_id INT
- fa_name VARCHAR(255)
- fa_archive_name VARCHAR(255)
- fa_storage_group VARBINARY(16)
- fa_storage_key VARBINARY(64)
- fa_deleted_user INT
- fa_deleted_timestamp BINARY(14)
- fa_deleted_reason TEXT
- fa_size INT
- fa_width INT
- fa_height INT
- fa_metadata MEDIUMBLOB
- fa_bits INT
- fa_media_type ENUM(...)
- fa_major_mime ENUM(...)
- fa_minor_mime VARBINARY(100)
- fa_description TINYBLOB
- fa_user INT
- fa_user_text VARCHAR(255)
- fa_timestamp BINARY(14)
- fa_deleted TINYINT
- fa_sha1 VARBINARY(32)
- Indexes

**uploadstash**
- us_id INT
- us_user INT
- us_key VARCHAR(255)
- us_orig_path VARCHAR(255)
- us_path VARCHAR(255)
- us_source_type VARCHAR(50)
- us_timestamp VARBINARY(14)
- us_status VARCHAR(50)
- us_chunk_inx INT
- us_props BLOB
- us_size INT
- us_sha1 VARCHAR(31)
- us_mime VARCHAR(255)
- us_media_type ENUM(...)
- us_image_width INT
- us_image_height INT
- us_image_bits SMALLINT
- Indexes

**oldimage**
- oi_name VARCHAR(255)
- oi_archive_name VARCHAR(255)
- oi_size INT
- oi_width INT
- oi_height INT
- oi_bits INT
- oi_description TINYBLOB
- oi_user INT
- oi_user_text VARCHAR(255)
- oi_timestamp BINARY(14)
- oi_metadata MEDIUMBLOB
- oi_media_type ENUM(...)
- oi_major_mime ENUM(...)
- oi_minor_mime VARBINARY(100)
- oi_deleted TINYINT
- oi_sha1 VARBINARY(32)
- Indexes

## Statistics

**hitcounter**
- hc_id INT

**site_stats**
- ss_row_id INT
- ss_total_views BIGINT
- ss_total_edits BIGINT
- ss_good_articles BIGINT
- ss_total_pages BIGINT
- ss_users BIGINT
- ss_active_users BIGINT
- ss_images INT
- Indexes

## Search

**searchindex**
- si_page INT
- si_title VARCHAR(255)
- si_text MEDIUMTEXT
- Indexes

## ResourceLoader

**module_deps**
- md_module VARBINARY(255)
- md_skin VARBINARY(32)
- md_deps MEDIUMBLOB
- Indexes

**msg_resource**
- mr_resource VARBINARY(255)
- mr_lang VARBINARY(32)
- mr_blob MEDIUMBLOB
- mr_timestamp BINARY(14)
- Indexes

**msg_resource_links**
- mrl_resource VARBINARY(255)
- mrl_message VARBINARY(255)
- Indexes

## Interwiki

**sites**
- site_id INT
- site_global_key VARBINARY(32)
- site_type VARBINARY(32)
- site_group VARBINARY(32)
- site_source VARBINARY(32)
- site_language VARBINARY(32)
- site_protocol VARBINARY(32)
- site_domain VARCHAR(255)
- site_data BLOB
- site_forward
- site_config BLOB
- Indexes

**site_identifiers**
- si_site INT
- si_type VARBINARY(32)
- si_key VARBINARY(32)
- Indexes

**interwiki**
- iw_prefix VARCHAR(32)
- iw_url BLOB
- iw_api BLOB
- iw_wikiid VARCHAR(64)
- iw_local
- iw_trans TINYINT
- Indexes

## Caching tables

**objectcache**
- keyname VARBINARY(255)
- value MEDIUMBLOB
- exptime DATETIME
- Indexes

**l10n_cache**
- lc_lang VARBINARY(32)
- lc_key VARCHAR(255)
- lc_value MEDIUMBLOB
- Indexes

**transcache**
- tc_url VARBINARY(255)
- tc_contents TEXT
- tc_time BINARY(14)
- Indexes

**querycache**
- qc_type VARBINARY(32)
- qc_value INT
- qc_namespace INT
- qc_title VARCHAR(255)
- Indexes

**querycache_info**
- qci_type VARBINARY(32)
- qci_timestamp BINARY(14)
- Indexes

**querycachetwo**
- qcc_type VARBINARY(32)
- qcc_value INT
- qcc_namespace INT
- qcc_namespacetwo INT
- qcc_title VARCHAR(255)
- qcc_titletwo VARCHAR(255)
- Indexes

## Maintenance

**updatelog**
- ul_key VARBINARY(255)
- ul_value BLOB
- Indexes

**job**
- job_id INT
- job_cmd VARBINARY(60)
- job_namespace INT
- job_title VARCHAR(255)
- job_timestamp VARBINARY(14)
- job_params BLOB
- job_random INT
- job_attempts INT
- job_token VARBINARY(32)
- job_token_timestamp VARBINARY(14)
- job_sha1 VARBINARY(32)
- Indexes

https://upload.wikimedia.org/wikipedia/commons/f/f7/MediaWiki_1.24.1_database_schema.svg

**revision**
- 🔑 rev_id INT
- ◈ rev_page INT
- ◈ rev_text_id INT
- ◈ rev_comment TINYBLOB
- ◈ rev_user INT
- ◈ rev_user_text VARCHAR(255)
- ◈ rev_timestamp BINARY(14)
- ◈ rev_minor_edit TINYINT
- ◇ rev_deleted TINYINT
- ◇ rev_len INT
- ◇ rev_parent_id INT
- ◇ rev_sha1 VARBINARY(32)
- ◇ rev_content_model VARBINARY(32)
- ◇ rev_content_format VARBINARY(64)

Indexes ▶

**page**
- 🔑 page_id INT
- ◈ page_namespace INT
- ◈ page_title VARCHAR(255)
- ◈ page_restrictions TINYBLOB
- ◈ page_counter BIGINT
- ◈ page_is_redirect TINYINT
- ◈ page_is_new TINYINT
- ◈ page_random DOUBLE
- ◈ page_touched BINARY(14)
- ◇ page_links_updated VARBINARY(14)
- ◈ page_latest INT
- ◈ page_len INT
- ◇ page_content_model VARBINARY(32)
- ◇ page_lang VARBINARY(35)

Indexes ▶

**page_props**
- ◈ pp_page INT
- ◈ pp_propname VARBINARY(60)
- ◈ pp_value BLOB
- ◇ pp_sortkey FLOAT

Indexes

**page_restrictions**
- 🔑 pr_id INT
- ◈ pr_page INT
- ◈ pr_type VARBINARY(60)
- ◈ pr_level VARBINARY(60)
- ◈ pr_cascade TINYINT
- ◇ pr_user INT
- ◇ pr_expiry VARBINARY(14)

Indexes

**protected_titles**
- ◈ pt_namespace INT
- ◈ pt_title VARCHAR(255)
- ◈ pt_user INT
- ◇ pt_reason TINYBLOB
- ◈ pt_timestamp BINARY(14)
- ◇ pt_expiry VARBINARY(14)
- ◇ pt_create_perm VARBINARY(6

**text**
- 🔑 old_id INT
- ◈ old_text MEDIUMBLOB
- ◈ old_flags TINYBLOB

Indexes ▶

**redirect**
- 🔑 rd_from INT
- ◈ rd_namespace INT
- ◈ rd_title VARCHAR(255)
- ◈ rd_interwiki VARCHAR(32)
- ◈ rd_fragment VARCHAR(255)

Indexes ▶

**category**
- 🔑 cat_id INT
- ◈ cat_title VARCHAR(255)
- ◈ cat_pages INT
- ◈ cat_subcats INT
- ◈ cat_files INT

Indexes ▶

https://upload.wikimedia.org/wikipedia/commons/f/f7/MediaWiki_1.24.1_database_schema.svg

nfldb Entity-Relationship diagram (condensed)

**game**
- gsis_id [gameid, not null]
- start_time [utctime, not null]
- week [usmallint, not null]
- season_year [usmallint, not null]
- season_type [season_phase, not null]
- finished [boolean, not null]
- home_team [varchar, not null]
- home_score [usmallint, not null]
- away_team [varchar, not null]
- away_score [usmallint, not null]

**drive**
- gsis_id [gameid, not null]
- drive_id [usmallint, not null]
- start_field [field_pos, null]
- start_time [game_time, not null]
- end_field [field_pos, null]
- end_time [game_time, not null]
- pos_team [varchar, not null]
- pos_time [pos_period, null]

**player**
- player_id [varchar, not null]
- full_name [varchar, null]
- team [varchar, not null]
- position [player_pos, not null]
- status [player_status, not null]

**team**
- team_id [varchar, not null]
- city [varchar, not null]
- name [varchar, not null]

**play_player**
- gsis_id [gameid, not null]
- drive_id [usmallint, not null]
- play_id [usmallint, not null]
- player_id [varchar, not null]
- team [varchar, not null]

**play**
- gsis_id [gameid, not null]
- drive_id [usmallint, not null]
- play_id [usmallint, not null]
- time [game_time, not null]
- pos_team [varchar, not null]
- yardline [field_pos, null]
- down [smallint, null]
- yards_to_go [smallint, null]

**meta**
- version [smallint, null]
- season_type [season_phase, null]
- season_year [usmallint, null]
- week [usmallint, null]

**agg_play**
- gsis_id [gameid, not null]
- drive_id [usmallint, not null]
- play_id [usmallint, not null]

home
away

0..N   0..N
0..N
1   1
1
1
1
1   1
0..N
0..N
1   0..N
1
0..N
0..N
0..N
1
0..N
1
0..N
0..N
1   1
1   1
0..N

https://github.com/BurntSushi/nfldb/wiki/The-data-model#er-diagrams

**Ownership**

```
{
  owners: [{
    type: SRE,
    user: jdoe
  }, ...]
}
```

**Dataset**

urn
platform
name
fabric

**OwnedBy**
type

1

n

**User**

urn
firstName
lastName

1

**HasAdmin**

1

**Group**

urn
name
size

n

**HasMember**

1

**Membership**

```
{
  admin: jdoe,
  members: [
    jdoe, ...]
}
```

**Profile**

```
{
  firstName: John,
  lastName: Doe
  ldap: jdoe
}
```

https://engineering.linkedin.com/blog/2019/data-hub

# All Variations of ER diagrams

In practice, everyone uses different notations.
What matters are the core *concepts*

(in this class, we will learn a specific notation)

# COURSEWORKS@COLUMBIA

**COMSW4111_001_2015_3: INTRODUCTION TO DATABASES (Fall 2015)**

View Site As  | ✓ – Select Role –
- Student
- Teaching Assistant

Home 🏠

Files & Resources 📁

Syllabus 📄

Mailtool ✉️

Gradebook 📒

Site Settings 🖼️

Library Reserves 📖

Research Guides 📖

Roster 📇

Textbooks 📘

Piazza ⬜

Help ❓

## ⇄ INTRODUCTION TO DATABASES

| Edit | Permissions |

**CourseNo:** COMSW4111_001_2015_3
**Meeting Time:** MW 02:40P-03:55P **Meeting Location:** SEELEY W. MU 833

**Instructor Information:**
Eugene Wu

COMSW4111_001_2015_3

# Entity-Relationship Modeling

Entities (objects) to store and their attributes

Relationships between entities and their attrs.

Integrity constraints & business rules

## ⇄ NEXT SEMESTER COURSES

**Fall 2015 – Spring 2016 Courses**

| Course Number | Course Title |
| --- | --- |
| COMSE6910_024_2015_3 | FIELDWORK |
| COMSW4111_001_2015_3 | INTRODUCTION TO DATABASES |

Reflects Registrar changes through Mar–06–2015 2:02:13AM

Courses
- Course Number
- Course Title
- Year
- Semester

**Eugene Wu** test test again just then [Clear](#)

| Say something | | **Say it** |
|---|---|---|

| **Profile** | **Wall** |
|---|---|

**Basic Information**

Nickname

Birthday

Personal summary

| **B** | *I* | <u>U</u> | ~~ABC~~ | X₂ | X² | | | | | HTML |

**Save changes**  **Cancel**

**Contact Information**

Email  ew2493@columbia.edu

Home page

Work phone

Home phone

Mobile phone

Facsimile

**Save changes**  **Cancel**

Users
  Nickname
  Name
  Birthday
  Summary
  Email
  …

# Basics: Entities

Entity e.g., intro to databases
    real-world object distinguishable from other objects
    described as set of attributes & the values
    (think one record)

Entity Set e.g., all courses
    collection of similar entities
    all entities have same attributes (unless Is-A)
    must have one or more keys
    attributes have domains
    ≈ table

# Example: Entity

Keys (cid, uid) are underlined

   Values must be unique

   (can use as hashtable key to lookup in table)

Course
  cid
  name
  loc
  schedule

Users
  uid
  name
  age
  summary

# Basics: Relationships

Relationship: association between 2 or more entities

    e.g., alice **is taking** Introduction to DBs

Relationship Set: collection of similar relationships

    N-ary relationship set R relates N entity sets $E_1 \ldots E_n$

    Each $r \in R$ involves entities $e_1 \ldots e_n$

    An $E_i$ can be part of diff. relationship sets or diff. roles in same set

# Basics: Relationships

Users can have different roles
in same relationship set

# Basics: Relationships

Relationships sets can have descriptive attributes

Denoted with dotted line from diamond to box

*since is attribute of Takes*

# Basics: Ternary Relationships

Connects three entities

N-ary relationships possible too.



*Assignments, Courses, and Users participate in the Graded relationship set*

# Constraints

Help avoid corruption, inconsistencies

Key constraints

Participation constraints

Weak entities

Overlap and covering constraints

# Key Constraints

Defines cardinality requirements on relationships

Many to many e.g., *Takes*

   a user can take many courses

   a course can have many users that take the course

One to Many e.g., *Instructs*

   a course has at most one instructor

*Draw arrow from diamond to box*



1-to-1          1-to Many          Many-to-Many

A course is instructed by ≤1 user
(read along the beige arrow)

A user instructs ≤1 course

A course is instructed by ≤1 user
AND
A user instructs ≤1 course

# Participation Constraints

Does every course need an instructor?

    If yes, it's a participation constraint

    Otherwise, *partial* participation constrain*t*

Denoted by double line between entity set and relationship set

*Each course MUST have an instructor*
*(participation of Courses in instructs is Total)*



*Each user must take at least one course and*
*Each course must have at least one user (student)*

# Weak Entities

A *weak entity* can only be uniquely identified by using the primary
key of its owner entity

    Owner and weak entity sets must have 1-to-N relationship

    Weak entity set must have total participation in this
    *identifying* relationships set

Denoted as double line around weak entity, set relationship set,
and the edge between them; an arrow to owner entity

# General Cardinality Constraints

Users ← Instructs — Courses

same as

Users —0:*— Instructs —0:1— Courses

A user instructs 0 to ∞ courses       A course has 0 to 1 instructors

A —x:y— ◇ —n:m— B

Each A entity has a relationship with between x to y different B entities
Each B entity has a relationship with between n to m different A entities

# Read arrows pointing in the direction from start to end

Each A is related to at most 1 B;  A has N-to-1 relationship with B

A  ⬦  B

Each B is related to any number of As;  B has 1-to-N relationship with A

A ⟵⬦— B    B has at most one A

A —⬦═ B    B has at least one A

A ⟵⬦═ B    B has exactly one A

A ⟵⬦═ B    B is a weak entity

# Specialization Hierarchies

Inheritance rules similar to programming languages

add descriptive attributes specific to a subclass e.g., grade

identify entity set that participate in a relationship

Denoted with arrow from subclass to superclass without a diamond



Every student & employee is a user

A user can be a student AND an employee

Every instructor & TA is an employee

An employee is *either* an instructor or TA

# Specialization Hierarchies

**Overlap Constraint**

can A be a B *and* a C?

YES           NO



*separate arrows*     *merged into 1 arrow*

**Total Specialization Constraint**

*must* A be a B or C?
specify as the comment "total"
with dashed link to arrows



total

# Aggregation

Relationships between (entities – relationships)

Treat Relationship Set like an Entity Set to participate in other relationships

Denoted as dashed line around the relationship set

# Aggregation vs Ternary Relationships

Why use aggregation?

Manages and Donates are distinct relationships with own attrs

Can define constraints on relationship sets

# Aggregation vs Ternary Relationships

Constraints apply to all connected entity sets

*A donation can be managed by at most one instructor*

But also enforces: *A course can have at most one donation*

# Using the ER Model

OK, we've seen the *syntax*.

How to use it involves design choices

Design Choices for a concept

    Entity or Attribute?

    Entity or Relationship?

    Binary or Ternary relationship?

    Aggregation or Ternary relationship?

# Entity or Attribute?

Is users.address an attribute of Users or an entity
connected to Users by a relationship?

Depends (and may change over time!)

If a user has >1 addresses, must be an entity

If an address has attrs (structure), must be entity

e.g., want to search for users by city, state, or zip

# Entity or Attribute?

A company can't donate multiple amounts

Company can make multiple donations

# Entity or Relationship?

But what if company donates to school for all data-related courses?

Redundancy of *amount*, need to remember to update every one

Misleading implies *amount* tied to *each* donation individually



| Company | Course | Amount |
|---------|--------|--------|
| Amazon | 4111 | 2000 |
| Amazon | 4112 | 2000 |
| Amazon | 5111 | 2000 |

*These amounts are logically the same (redundant)!*

# Entity or Relationship?

If company donates once to school for data related courses.
Refactor amount into an entity



Companies — Donates — Courses

Donations
Amount
When

*Company redundant, since company only donated once*

| Company | Course | Donation |
|---------|--------|----------|
| Amazon | 4111 | 1 |
| Amazon | 4112 | 1 |
| Amazon | 5111 | 1 |

| Donation | When | Amount |
|----------|------|--------|
| 1 | Today | 2000 |

# Entity or Relationship?

If company donates once to school for data related courses.

Refactor amount into an entity



| Course | Donation |
|--------|----------|
| 4111   | 1        |
| 4112   | 1        |
| 5111   | 1        |

| Donation | When  | Amount | Company |
|----------|-------|--------|---------|
| 1        | Today | 2000   | Amazon  |

# Binary or Ternary Relationship?

What if each HW has at most one grader?

# Binary or Ternary Relationship?

What if each HW has at most one grader?



*Actually says that each student's HW submission has at most one grader*

*Each HW has at most 1 grader and the grader evaluates student submissions*

# Binary or Ternary Relationship?

Binary relationships allows additional constraints

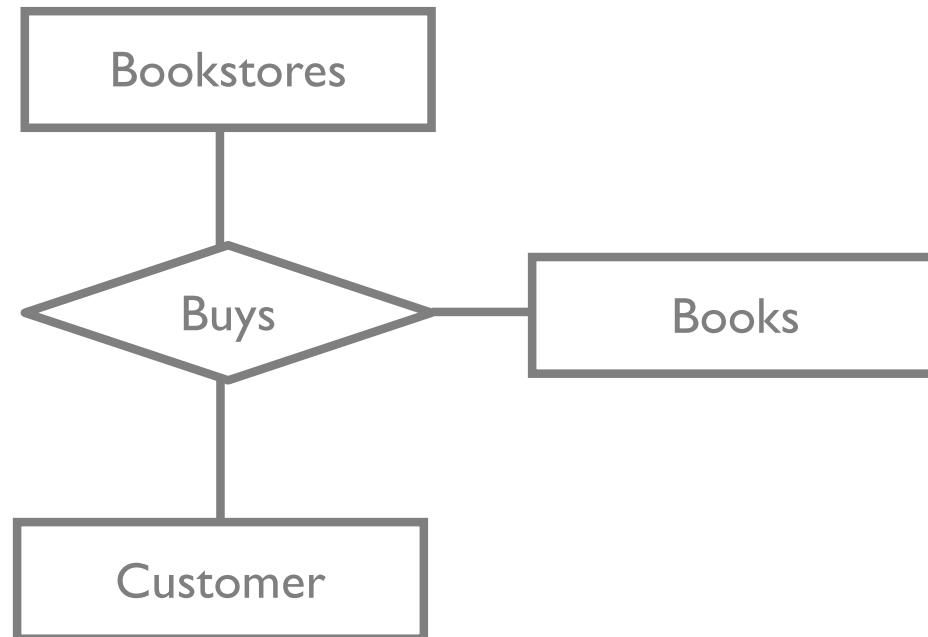What should happen if a student drops the class? (see next slide)

# Binary or Ternary Relationship?

Binary relationships allows additional constraints



*When student drops the class, HW0 also disappears! Previous slide was correct*
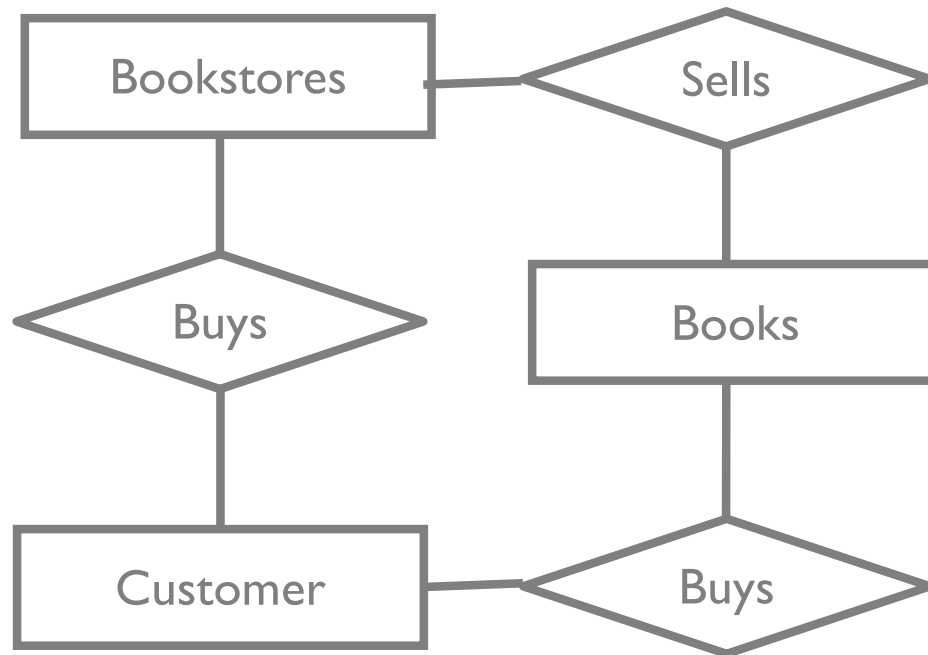
# Binary or Ternary Relationship?

Sometimes have true ternary relationship that is defined by all three entities.

# Binary or Ternary Relationship?

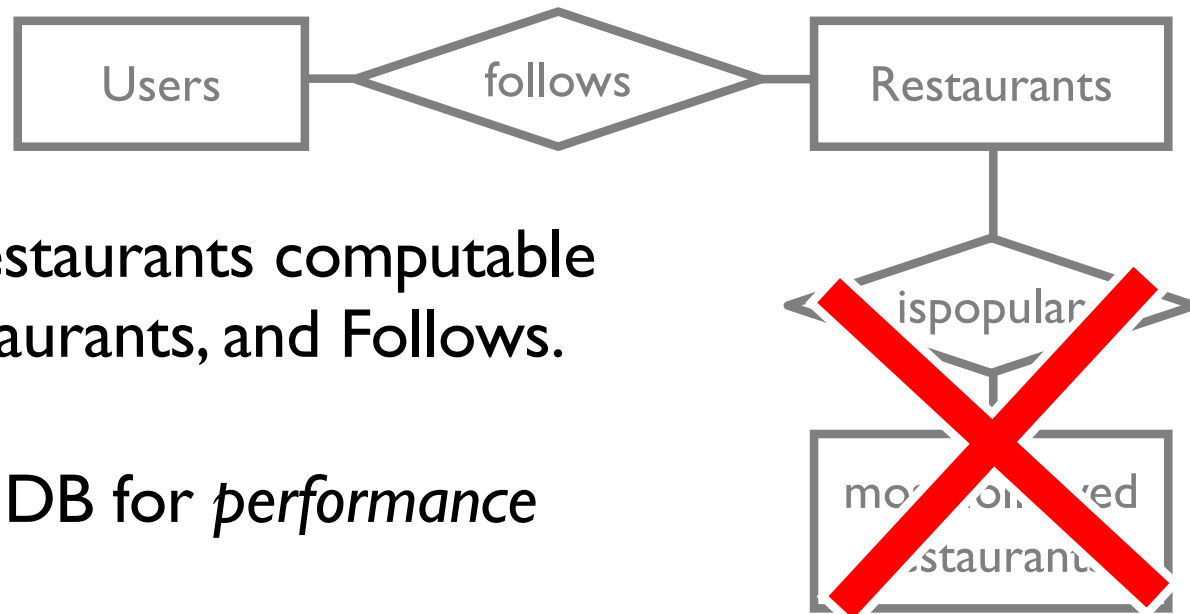Sometimes have true ternary relationship that is defined by all three entities.

*Doesn't Really Work*

# Advice

The ER diagram (and database) stores the *minimal information* needed for your application.

Everything else (e.g., stats) can be computed



Most followed restaurants computable from Users, Restaurants, and Follows.

May still store in DB for *performance* reasons

# Summary

Requirements
    what are you going to build?
Conceptual Database Design          (Today) ER Modeling
    pen-and-pencil description
Logical Design
    formal database schema
Schema Refinement:
    fix potential problems, normalization
Physical Database Design
    use sample of queries to optimize for speed/storage
App/Security Design
    prevent security problems

# Summary

Conceptual design follows *requirements analysis*

ER model helpful for conceptual design

    constraints are expressive

    matches how we often think about applications

Core constructs

    entity, relationship, attribute

    weak entities, ISA, aggregation

Many variations beyond today's discussion

# Summary

ER design is subjective based on usage+needs

    Today we saw multiple ways to model same idea

ER design is not complete/perfect

    Developed in an enterprise-oriented world (ER First)

    Doesn't capture semantics (what does "instructor" *mean?*)

    Doesn't capture e.g., processes/state machines

    How to combine multiple ER models automatically?

    Limitation of imagination when designing application

    Still needs further refinement

    Open problems!

ER design is a useful way to think

# Next Time

Relational Model: de-facto DBMS standard

Set up for ER diagrams → Relational models