# Information Architecture and Search

INFO 200

Part I

**Joseph Janes**
**Associate Professor, Information School**

---

# Information Architecture & Search 1
agenda

⬧ search, and what makes it work:  structured search
  ⬧ data modeling & encoding
  ⬧ database management systems

---

# information behaviors

information **use, seeking, retrieval, organization, encountering, etc**

all imply **seeking** or **finding** at some level  (so does information **destruction, censorship** for that matter)

- so **what enables search?** what makes search work, makes it possible?

(and, by the way, all of what is to come are information behaviors too)

---

# search, and what makes it work

for what?
  *information objects*
fair enough - what are information objects?
  *Web pages, books & published materials, tweets*
  *words, facts*
  *people, organizations*
  *sounds, images, moving images, objects*
  *basically anything*
different objects, different searches
the search for each of these is enabled differently, with common features
an example:  General Education Requirement Course Search

# search

**search** is the **matching** of **representations** in a **database** by means
  of an **algorithm**

(representations = **metadata**)

this implies those representations, databases, algorithms have been
  somehow **created, defined, decided on** (all information behaviors)
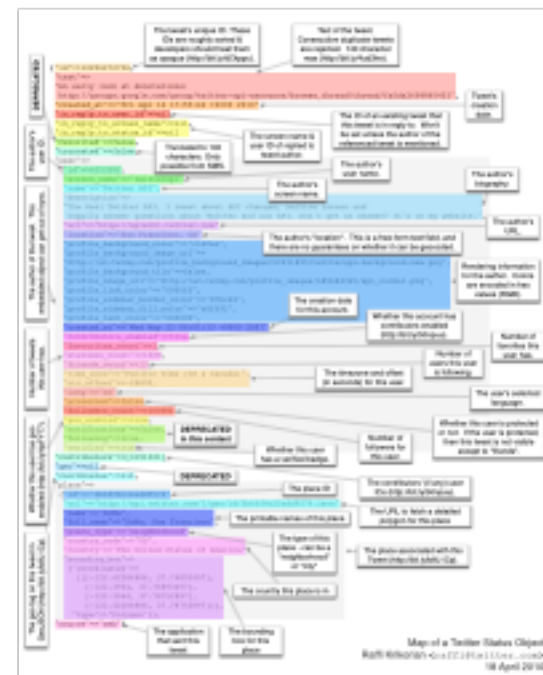
and **structured**

so this "simple" instance is referred to as **structured search**

more examples:

  UW Faculty/Staff/Student Directory - Seattle

  Olympedia

**no metadata, no search**

these are all examples of **information systems** - so where is **power**?

## Jacquard Loom
1804



## Hollerith punch card
1889

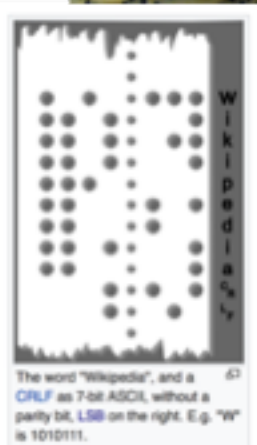## US Census schedule
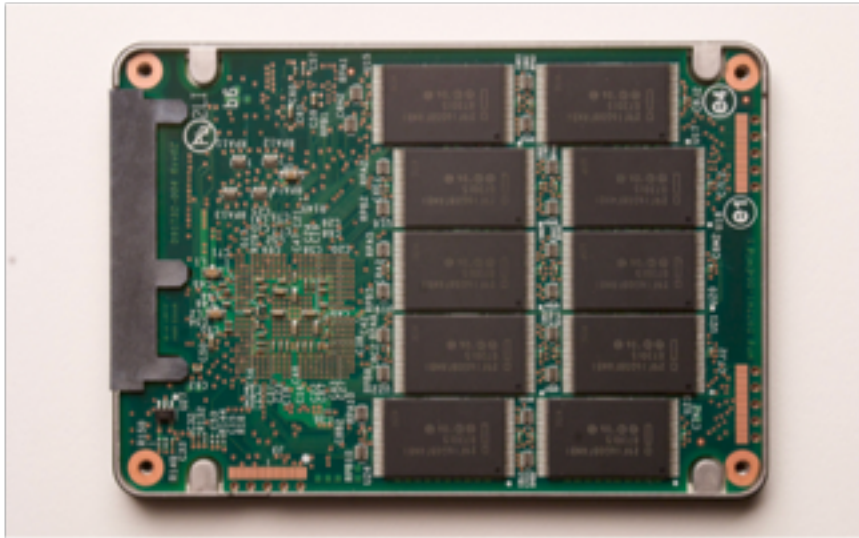### 1940



## IBM punch card
### 1928



## punch paper tape
### orig. 1725 looms; computer/communication by 1944



The word "Wikipedia", and a CRLF as 7-bit ASCII, without a parity bit, LSB on the right. E.g. "W" is 1010111.

If computers store only patterns of bits, how do we reliably encode **text** into files so that multiple programs can display that text again?

# character sets

Define a mapping between patterns of bits and characters

Contain decisions that may have significant social ramifications



ASCII Character Set (1963)
7 bits per character = 128 possible characters

# multiple character sets cause serious problems

| bits | encoding | characters |
|---|---|---|
| 11000100 01000010 | Windows Latin 1 | ÄB |
| 11000100 01000010 | Mac Roman | ſB |
| 11000100 01000010 | GB18030 | 䏮 |

| characters | encoding | bits |
|---|---|---|
| Foō | Windows Latin 1 | 01000110 11111000 11110110 |
| Foō | Mac Roman | 01000110 10111111 10011010 |
| Foō | UTF-8 | 01000110 11000011 10111000 11000011 10110110 |

## Unicode (1991)

16 bits per character

65,536 possible "code points" (characters)



http://unicode-table.com

## multiplane Unicode



32 bits per-character = 4.2 billion code points

Multiple encoding strategies:

| | |
|---|---|
| UTF-32 | 32 bits for every character (UNIX) |
| UTF-16 | 16 bits for low chars; 32 for high (Java, .Net) |
| UTF-8 | 8 to 32 bits, depending on char (the web, Go) |

## unicode encoding strategies

| character | encoding | bits |
|---|---|---|
| A | UTF-8 | 01000001 |
| A | UTF-16 | 00000000 01000001 |
| A | UTF-32 | 00000000 00000000 00000000 01000001 |
| あ | UTF-8 | 11100011 10000001 10000010 |
| あ | UTF-16 | 00110000 01000010 |
| あ | UTF-32 | 00000000 00000000 00110000 01000010 |

## raster image encoding



Red: 233
Green: 157
Blue: 144

3 numbers per pixel, each 0-255

24 bits per pixel

# sound encoding



Bit depth is the number of **bits** available for each **sample**. The higher the bit depth, the higher the quality of the audio. Bit depth is usually 16 bits on a CD and 24 bits on a DVD.

A bit depth of 16 has a resolution of 65,536 possible values (ranging from 0 to 65,535), and a bit depth of 24 has over 16 million possible values (ranging from 0 to 16,777, 216).

The bit rate is calculated using the formula:

**Frequency × bit depth × channels = bit rate**

A typical, uncompressed high-quality audio file has a **sample rate** of 44,100 **samples** per second, a bit depth of 16 bits per sample and 2 channels of stereo audio. The bit rate for this file would be:

**44,100 samples per second × 16 bits per sample × 2 channels = 1,411,200 bits per second (or 1,411.2 kbps)**

A four-minute (240 second) song at this bit rate would create a file size of:

**1,411,200 × 240 = 338,688,000 bits (or 40.37 megabytes)**

https://electronics.howstuffworks.com/analog-digital3.htm
https://www.bbc.com/bitesize/guides/z7vc7ty/revision/1

ⁿ`¬«-ⁿΩⁿAᵖ⁴ⁱⁿⁱⁿⁿᵗⁿⁿqⁿⁿⁿⁿbⁿⁿxⁿ¶
¾z¿½É<-/10°ør⌐¾ä Ûûû-}2Ï9='⌐é3ⁿⁿⁿⁿÀ•6ᵍ,ⁿⁿ⌐ⁿⁿ436>Ê‡Àⁿⁿⁿ4ⁿ(ⁿ-ÛⁿⁿÛⁿ+XⁿCⁿⁿⁿⁿⁿÄ=",¸'ⁿ¸µⁿ«ⁿ¿'ⁱ˜¢ⁿÉⁿ¾;
¸ÔⁿⁿÛⁿⁿⁿⁿⁱ'ⁿ88ⁿⁿⁿ¨,Aⁿⁿⁿⁿ=Xⁱ=(ⁱⁿ,ⁱÛⁿ†Ûⁿⁿ);ⁿ¼)[ⁿÄû)Éⁿⁿⁿz(ⁿ(ⁿⁿ)ⁿⁱ'Sⁿⁿⁿⁿⁿⁿⁿ'⌐ⁿÛ‰ⁿⁿⁿⁿ'ⁿ‴ⁿⁿⁿⁿⁿⁿⁿ2 =ⁿⁿⁿⁿⁿⁿⁿ+†ⁿÛⁿ'ⁿⁿ:¸ⁿÉⁿ§sⁿⁿ •ⁿⁿⁿⁿⁱⁿ
}¸-⌐ⁿⁿⁿÛⁿⁿ.
ƒ//(Ûⁿ-†˜ⁿⁿⁿ,ƒⁿⁿⁿⁿ=
¸ⁿⁿⁿⁿⁿ'Û=ÇÁÛⁿⁿⁿⁿⁿ,˜¶ⁿⁿⁿⁱⁿ/ⁿⁿⁿⁿⁿⁱ
ÂⁿÛⁿⁱ>±ⁱⁿⁿⁿ7ⁿⁿⁿ±ⁿⁱⁿⁿⁿⁿⁿ⌐ⁿ¶ⁿⁿⁿ‡ⁿ«=-
=ⁿⁿÄr=ⁿⁿⁿⁿⁱⁿⁿ(ⁿ¾ⁿ‡ⁿÛⁿⁿ¹ⁿⁿ<-=²•ⁿⁿ*ⁿⁿⁿⁿⁿⁿ;J¸ⁿ±ⁱⁿⁿⁿⁿⁿ⌐ⁿⁿⁿ˜•ⁿ'Éⁿ0
ⁿⁿⁿ ⁿ˜ⁿⁿⁿⁿⁱⁿ-'ⁿⁿⁿ>ⁿⁱ,ⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿⁱⁱ-Âⁿ¸ⁿⁱ,ⁿⁿⁿⁿⁿ,       ⁿⁿ=ⁿⁿ Iⁿⁿ‡ⁿⁿ)<ⁿ˜ⁿⁿⁿⁿ0ⁿⁿ>ⁿ;ⁿⁿⁿⁿⁿ ⁿ-¸>'ⁿⁿ
ⁿⁿⁿⁿⁿⁿÛ,Û1R¸ⁿ2ⁿⁿⁿ'ⁿ'ⁱⁿ⌐ⁿ¸<-ⁿ†ⁿ†ⁿⁿ0ⁿ;'--ⁿ¼Â=ⁿⁿⁿ'}ⁱⁿⁿ¸Ûⁿⁿ)ⁿⁿ=ⁿⁿ'ⁿ
ⁿⁿⁿ<=¸ⁿ0ⁿ⁰'†1 Â-(˜-ⁿⁿ7ⁿⁿⁿⁿ/ⁿ¢ⁿ†˜ⁿⁿ,ⁿ¸ⁿⁿⁿ=ⁿⁿⁿⁿⁿⁿ†⌐[ⁿⁿⁿⁿⁿⁿ'ⁱⁿⁿ¢ⁿ  =ⁿⁿⁿ=ⁿ‡ⁱⁿ Iⁿⁿⁿⁿⁱⁿ
1TⁿⁿⁿⁿÂ-ⁱⁿ1
¸ⁿⁿⁿⁿLⁿⁿ'ⁿ Qⁿⁿⁿ=ⁿⁿⁿⁿYⁿⁿⁿⁿⁿⁿⁿⁿⁿ¶ⁿⁿⁿ<ⁿ Δ¸'ⁿ'Éⁿⁿ‡';ⁿ¶ⁿⁱ+'\ⁱ?*'"¸ⁿⁿ'Ûⁿⁿⁿⁿⁿ='¸¸ⁿⁿⁱ'-ⁿ{\ⁱⁿ/¿:Hⁿⁿ,:rⁿⁿ)ⁱ=Êⁿⁿ⌐ⁱⁿ{ƒ
ⁿÄ¸ⁱ'ⁱ±ⁿⁿ1ⁱⁿ}ⁱⁿⁿ2ⁿⁿⁿⁿⁿⁿÛ-¸ⁿⁿⁱ¿ⁿ¶ⁿ†ⁱ=7ⁱⁿ¢¸=ⁿⁿⁿⁿⁿ'=-=ⁿⁿ¸ⁿⁿⁿⁿ-ⁿ¸Ûⁿⁿ(ⁿⁿⁿⁿⁿⁿⁿ'ⁿⁿⁿ¸',ⁿÉⁿ
'CⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿC}-¸ⁿⁿ ˜' =Xⁱ'ⁿ¶Sⁿⁿⁿⁿ¶ⁿⁿⁿⁿⁿ,"ⁿⁿÛ=ⁿÉⁿⁿⁿⁿ7 ˜ⁿⁿⁿⁿⁿ'ⁱⁿ'-
Äⁿⁿ'ⁿⁿ33ÛⁿⁿⁿⁿⁿÑⁿⁿⁿⁱⁿ¸'ⁿⁿ Ûⁿⁿⁿⁿⁿⁿ=ⁿ‡ⁿⁿ,\ⁿⁿⁿ-ⁿⁿⁿⁿ‡ⁿ¸Éⁿⁿⁿ[/¶ⁿⁿⁿ=ⁿⁿ='ⁿⁿⁿⁿⁿⁿⁿ5ⁿ[Âⁿ†ⁿⁿⁿⁿⁿⁿ4]ⁱⁿ/¸ⁿⁿⁿ¸ⁱ''˜ⁿⁿⁿxⁿⁿ"=ⁿⁿ'<Ûⁿ,¸ⁿⁿⁱ'ⁿ=Aⁿⁿrⁿⁿ-ÛⁿⁿPⁿⁿⁿⁿⁿⁿ'-(ⁿⁿÛÄ
<ⁿ7˜ⁿ(ⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿ ⁿⁿ†ⁿⁿⁿⁿⁿⁿ¸ⁿⁿ2˜ⁿⁿⁿⁿⁿ1ⁿⁿⁿⁿⁿⁿⁿ=
ⁿⁿⁱⁿ‡ⁿⁿⁿ¸-ⁱⁿⁿⁿⁿⁿⁿ¸'Ûⁿ¸ⁿⁿⁿⁿⁿⁿⁿⁿ'(ⁿ+ⁿ¸ⁿⁿⁿ<¸,ⁿⁿⁿⁿⁿⁿ˜ⁿⁿⁿⁿⁿ]"')Ûⁿ'ⁿⁿⁿ¸Aⁿⁿⁿ
ⁿⁿⁿⁿ=ⁿⁿⁿ,Ûⁿⁿⁿ?
ⁿⁿ{ⁱⁿⁿⁱⁿ,/ⁿⁿ'Rⁿⁿ¸˜"ⁿÛ]ⁿ'Anⁿⁿⁿ='ⁿ¸Zⁿⁿⁿⁿⁿⁿⁿ'ⁿ¸¸ⁿ'Tⁿⁿⁿⁿ'0'/ⁿ¶ⁿⁿⁿ'<ⁿⁿⁿ‡'
ⁱⁿ¸'ⁿⁿÉⁿⁿ Aˣⁿⁿⁿ/'(7/ⁿⁿⁿ=Fⁿⁿⁿⁿⁿ"Cˊ)˜ⁿ¸;"ⁿⁿⁿÛⁿ
ⁿⁿ-¸ⁿⁿⁿⁿ-˜ⁱⁿⁿⁿ"ⁱⁿⁿⁿ"{ⁿⁿⁿⁿ˜ⁿⁿⁿⁿⁿⁿ/ Äⁿⁿˊⁿⁿⁿⁿ[ⁿ'ⁿ't˜ⁿ?Zⁿⁿⁿ§†ⁿⁱ[Ûⁿ§ⁿ§ⁿⁿⁿ)/1Gˊⁿⁿ-û",ⁿˊⁿƒ/ⁿ(ⁿ«ⁿⁿⁿⁿⁿⁿ=ⁿ
Ⅴ¸ˊⁿⁿⁿÈ['¸2¥ⁿⁿ°ⁿˉ˜N"-ⁿ
Ⅾ1ⁿⁿ%=</ⁿ1Ûⁿⁿ,ⁱˉ'ma}1fⁿ¹,<Éˉ¸ⁱⁿ˜';rⁿⁿⁿⁿ N˜Aⁿ§˜-¾4
ⁿⁿⁿⁿⁿ'PⁿÛ˜V"ⁿⁿⁿ¸"'-ⁿ}ⁿ,Ç
²«f'ⁿⁿ///ⁿⁿ¶ⁱⁿ†ⁿ'ⁿⁿⁿⁿ7Ⅾ}ⁿ1¸'"Ⅴ'ⁿⁿⁱⁿⁿⁿⁿⁿ=ⁿ[ⁿ†ⁿ¶ⁿ\ⁿ§ⁿ=ⁿⁿⁿⁿ-<ⁿ#ⁿ"ⁿ#†ⁿⁱⁿ=ⁿⁿˣ;J¸-Äⁿ¸"1]§ⁿ¸ⁿⁿÉ‡ⁿ>µⁿⁿ'Äⁿⁿ(X§rˉ'pⁿ7
Zⁿⁿ¸A ¢-ⁱ,S1Qⁿⁿⁿ«ⁱⁿ¸ⁿⁿⁿⁿ4ⁿ¸9ⁿⁿ=É((mⁿⁿ='SⁿⁿⁿⁿⁿⁿⁿⁿⁿÉⁿ"]ⁿZ¸Zⁿⁿⁿ˜=ⁿ'ⁿⁿ¾5ⁿ:ˉⁿⁿⁿ=2ⁿ‡ⁿⁿ=ⁿⁿⁿⁿ7ⁿⁿ¸¸ⁿ=ⁿⁿ§ⁿⁿⁱ'(ⁿⁿⁿ$
¸ⁿⁿ=ⁿⁿⁿⁿⁿÀⁿⁿⁿⁿⁿ-Ûⁿⁱ$ⁿⁿˉ;ⁿⁱⁿ}˜˜ⁿⁿⁿ$ⁿ/ⁿⁿ¸Â¸Àⁿ¸¸-˜'ⁱ-ⁿ2ⁿ†ⁱⁿⁿÉ¿QⁿI=ⁱⁱⁿ"ƒƒ¢ˊⁿ,Iⁿⁿⁱⁿⁿⁱ}ⁿⁿⁿÛⁿ'ⁿⁿⁿG)ⁱ˜'Ûˊⁿⁿⁿⁿˊⁱⁿⁿ\ˊⁿ¸ⁿⁱⁿⁿ'ⁿ
˜ⁿⁿⁿⁿⁿⁿⁿⁿÄⁿ /ⁿⁿⁱⁿ=ⁱⁿⁿⁿⁿⁿⁿⁱⁿⁿⁿˊⁱ[ⁿⁿ+ⁿ,mⁿ{ⁱⁱ=ⁿˊⁿ¸ⁿ⁰>U˜
µⁿÉˉ-ⁿⁿⁿ,=ⁿⁱⁿˊ3¸ⁿ•É}ⁿⁿ,ⁿⁿ,sⁿˉCˉ5=ⁿ¸XⁿⁿˉJⁿⁿ
ⁿⁿ]¸ⁿⁿⁿⁿˉ
ⁿⁿ{ⁿ=ˉⁿⁿⁱV¸Yⁿⁿⁿ ⁱˉ:
ⁿⁿⁿⁿⁿ,Ç'ⁿ0ⁿⁿ=ⁱⁿⁿⁿ "ⁿ'Sⁿⁿⁿⁿⁿⁿⁿⁱⁱⁿ-ⁿ¢'ƒ]ⁿ¶}ⁿ{ⁿⁿⁿⁿˉⁿⁿⁿ
ⁿⁿⁿⁿ-ⁱⁿⁱⁿÀⁿ,-ⁿⁿ¸ⁿ‡ⁿⁿⁿ=§˜=ⁿⁿⁿ/ⁿⁿⁿ/ⁿⁿⁿⁱⁿ‡ⁿ!:--ⁿⁿ'ⁿⁿⁿ§ⁿⁿ"ÂⁿⁿⁿⁿⁿLⁿⁿⁿⁿ§ⁿⁿ7ⁿ'µAⁿⁿⁿⁿⁿⁿⁿⁿⁿⁿ†ⁿÛˊ'ⁿⁿⁿ,ⁿ-ⁿⁱⁿⁱⁿⁿⁿ7ⁿⁿⁿˉ-ⁿ'Cⁿⁱ7ⁿˉ-ˉⁿⁿⁿⁿ=ⁿⁿⁿⁿ=--ⁿⁿⁿ/ⁿⁿⁿⁿ73}-
Ûⁿⁿⁿⁿⁿⁿⁿⁱⁿⁿ
Ûⁿⁿⁿⁿⁱ/ⁿⁿⁿⁱⁿⁿ1Û=ⁿⁿ§ⁿⁿⁿᵍT,"=ⁿⁿ¸ⁿⁿⁿ˜'cⁿⁿⁿⁿⁿⁿ.ⁿ"2/¸ˉⁱⁿⁿ:4ⁿⁿⁿ3Ⅰⁿⁿˊ‡÷ˉⁿⁿ
ⁿⁿⁿⁿ<Û'ˉÛ'Z¸'Ûⁿ'ⁿⁱ<ⁿ=ⁿⁿⁿⁿⁿⁿⁿⁿ=ⁿⁿ'ⁿⁿ=—=ⁱⁿⁿⁿ}ⁿ˜Ûⁱⁿ' '‡ⁿⁱⁿⁿ,7ⁿ7/Ⅰⁿˊ4•Éⁿⁿ #Éⁿⁿⁿⁿⁱ˜
ⁿⁿⁿⁿⁿ}5)ⁿ3ⁿⁿⁿⁿ}ⁿ¸ⁿ§Ⅱⁿ‡ⁿⁿⁿⁿⁿJ-ⁿⁿ'ⁿˉ'ⁿ¸ⁱ!ⁱⁿⁿⁿⁿⁿⁿⁿ¸'ⁿⁿ
ⁿⁿⁿ>ⁿ«¸ⁿⁿⁿⁿ7ƒ
ⁿⁿⁿⁿ˜'ⁿⁿⁿⁱ>ⁱ¸

# representation

anything stored, digitally or otherwise, has to be **represented** in a structured, standardized way - so that it can be accessed and found

("representation" has multiple meanings, yes?)

we've seen representing text, sound, images, moving images - now another kind at a higher conceptual level

# Database Management System (DBMS)

A software process that allows clients to define a data structure; add data to that structure; navigate, tabulate, update, and delete those data; maintain data security and integrity; and automatically recover from failure

| data file | data file | data file | data file | data file |
|-----------|-----------|-----------|-----------|-----------|

**DBMS**

⇧ ⇧ ⇧

| Client 1 | Client 2 | Client 3 |
|----------|----------|----------|

---

# relational data modeling

The process of designing a structure capable of holding the data the system needs to track, while avoiding redundancy and ensuring data integrity

Resulting structure is called a **schema**

Visually designed and represented in an **Entity-Relationship Diagram (ERD)**
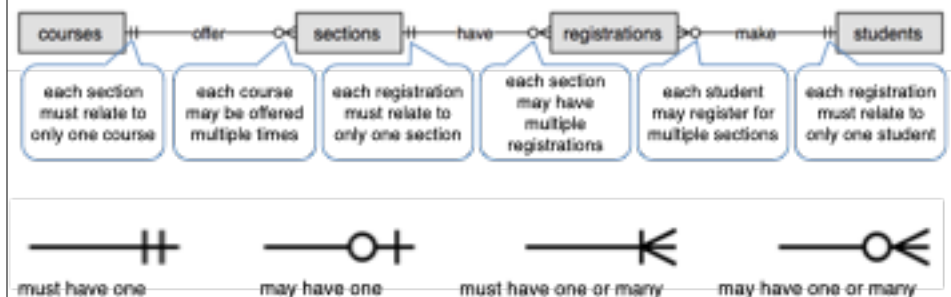
---

# entity

A person, place, thing, or concept included in a system

Look for the **core concepts** as people talk about their data, as well relationships that have their own data

"I want to track the courses we can teach, the sections of those courses we offer, the students who register for those sections, and what their final grades were."
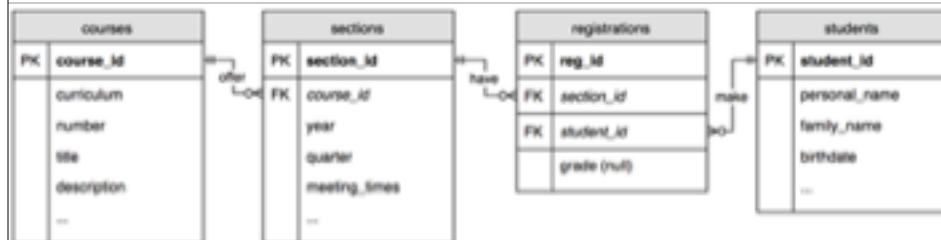
---

# relationships

How do entities relate, and how many instances can they relate to?

| courses | ⊦⊦ | offer | o< | sections | ⊦⊦ | have | o< | registrations | ⊳o | make | ⊦⊦ | students |
|---------|----|-------|----|----------|----|------|----|--------------|----|------|----|---------|

| each section must relate to only one course | each course may be offered multiple times | each registration must relate to only one section | each section may have multiple registrations | each student may register for multiple sections | each registration must relate to only one student |
|---|---|---|---|---|---|

⊦⊦ must have one     —o+ may have one     —⟨ must have one or many     —o⟨ may have one or many

## attributes

What do you want to track about each entity?



PK = Primary Key = Unique record identifier
FK = Foreign Key = PK value in related table
**other examples of keys?**

attributes with (null) allow missing data (null values)

---

## representation

anything stored, digitally or otherwise, has to be **represented** in a structured, standardized way - so that it can be accessed and found ("representation" has multiple meanings, yes?)

---

# US Census schedule
1940



---

# 1940 US census instructions



**INSTRUCTIONS FOR FILLING OUT THE POPULATION SCHEDULE**

*General Instructions*

24. Use *black* ink. Write legibly and keep your schedules neat and clean. Make all entries carefully.
25. Study very carefully the headings of all questions on the schedule and the symbols and explanatory notes at the bottom of the schedule.

6. **The Census Day.**—There should be a return on the Population Schedule for each person alive at the beginning of the census day, i. e., 12:01 A. M. on April 1, 1940.
7. **Who is to be enumerated.**—Enumerate all men, women, and children (including infants) whose usual place of residence (the place where they "*live*" or have their "*home*") is in your district, including persons temporarily absent; all persons who are in your district at the time of the enumeration who have no usual place of residence elsewhere from which they will be reported; and all persons who move into your district after the enumeration begins and who have not previously been enumerated. Enumerate as residents of the institution all inmates of a jail, however short their term of sentence, and all inmates of a prison, home for orphans, or similar institution located *in your district* in which persons remain for long periods of time.

https://1940census.archives.gov/

# 1940 US census instructions

> 40. **How names are to be written.**—Enter the last name first, then the given name and initial, making sure that the spelling is correct. Where the surname is the same as that of a member of the same household enumerated on a line above, do not repeat the name but enter a long dash. (See Illustrative Example, Form P-2.) Where there are not enough lines left on a schedule to enumerate all members of the household, fill out that side of the schedule completely, leaving no line vacant, make a check in the box labeled "Household continued on next page" in the lower left-hand margin of the schedule, and write "*Contd.*" (for "Continued") in cols. 1 and 2, (and leave cols. 3 to 6 blank) at the top of the "*B*" side of the schedule, or at the top of the
>
> 41.
>
> *Personal Description*
>
> 44. **Column 9. Sex.**—Write "*M*" for male and "*F*" for female.
>
> 45. **Column 10. Color or race.**—For symbols to be entered in this column, see the note at the bottom of the schedule. Any mixtures of white and nonwhite blood should be recorded according to the race of the nonwhite parent. A person of mixed Negro and Indian blood should be reported as Negro unless the Indian blood greatly predominates and he is universally accepted in the community as an *Indian*. Other mixtures of nonwhite parentage should be reported according to the race of the father. Mexicans are to be returned as *white*, unless definitely of Indian or other nonwhite race.

https://1940census.archives.gov/

# Information Architecture and Search

INFO 200

Part I

## Joseph Janes
### Associate Professor, Information School