

# DATA INTEGRATION

## THEORY REVIEW

THE AIM OF DATA INTEGRATION SYSTEMS (DISS) IS TO SET UP A SYSTEM WHERE IT'S POSSIBLE TO QUERY DIFFERENT DATA SOURCES AS THEY WERE A UNIQUE ONE THROUGH A GLOBAL SYSTEM

- ① HOMOGENEOUS DATA SOURCES: ALL SOURCES USE THE SAME DATA MODEL
- ② HETEROGENEOUS DATA SOURCES: DIFFERENT DATA MODELS

## MAIN STEPS

- ① SOURCE SCHEMA REVERSE ENGINEERING (LOGICAL → CONCEPTUAL SCHEMA)
- ② SCHEMA INTEGRATION (RELATED CONCEPT IDENTIFICATION + CONFLICT ANALYSIS)
- ③ GLOBAL CONCEPTUAL SCHEMA
- ④ CONCEPTUAL → LOGICAL OF THE GLOBAL SCHEMA
- ⑤ (GAV) MAPPING DEFINITION AND QUERY ANSWERING
  - a) GAV MAPPING DEFINITION
  - b) QUERY FORMULATION (ON GLOBAL SCHEMA)
  - c) QUERY REWRITING
- ⑥ CONFLICTS (NAME, TYPE, DATA SEMANTIC, STRUCTURE, CARDINALITY, KEY)

# Technologies for Information Systems

## Part II

prof. L. Tanca – February 28, 2012

Available Time 2h

Last Name	<hr/>
First Name	<hr/>
Student ID	<hr/> Signature

*LALuxuryHouses* is a real estate agency located in Los Angeles and its business is **exclusively** focused on luxury villas located in the Los Angeles area (State of California). Differently, *USA Houses* is an important real estate agency that rents and sells houses in all the main states of the USA. *USA Houses* wants to increase its business in Los Angeles. Since the Los Angeles area is currently only partially covered by the agencies of *USA Houses*, its management decided to buy *LALuxuryHouses* and founded a new company called *USARealEstateCompany*. The management of *USARealEstateCompany* (the new company) wants to integrate the information available in the two sources (*LALuxuryHouses* and *USA Houses*) in order to be able to query all the available data.

In the following we report the original relational schemas of the two sources.

### *LALuxuryHouses:*

CLIENTS (SSN, Lastname, Firstname, Address, City, State, Age, PhoneNumber)

EMPLOYEE (IDEmployee, Lastname, Firstname, PhoneNumber)

HOUSES (HouseAddress, HouseCity, SizeSquareMeters, Rooms) // *The size of each home is measured in square meters.*

HOUSE-OWNEDBY (HouseAddress, HouseCity, ClientSSN) // *Table House-OwnedBy is used to store the information about the owners of each house.*

RENTAL-CONTRACT (IDRentContract, HouseAddress, HouseCity, StartDate, EndDate, AnnualCost, IDEmployee) // *Each tuple in Table Rental-Contract represents the rental of a house (identified by the pair HouseAddress, HouseCity) for the period from StartDate to EndDate*

RENTEDBY (IDRentContract, ClientSSN) // *Table RentedBy is used to store who are the clients associated to each rental contract (i.e., who rented the house associated to the contract).*

SALE (IDSaleContract, HouseAddress, HouseCity, Date, Cost, IDEmployee) // *Each tuple in Sale corresponds to one sale.*

SOLDTo (IDSaleContract, ClientSSN) // *Table SoldTo is used to store who are the buyers associated to each sale.*

### **USA Houses:**

BUYERS (BuyerID, Name, Surname, Address, City, State, YearOfBirth, SSN, PhoneNumber)  
// Each tuple in Table Buyers represents someone who bought or rented a real estate  
OWNERS (OwnerID, Name, Surname, Address, City, State, YearOfBirth, SSN, PhoneNumber)  
// Each tuple in Table Owners represents someone who owns a real estate  
AGENTS (AgentID, Name, Surname, MobilePhoneNumber, OfficePhoneNumber)  
**REALESTATES** (IDRE, Address, City, State, NumOfRooms, Size\_SquareFeet,  
NumberOfFloors, OwnerID) // *The size of each real estate is measured in square feet.*  
REALESTATE-RENTAL (IDRE, StartDate, EndDate, BuyerID, AgentID, MonthlyCost)  
REALESTATE-SALE (IDRE, Date, BuyerID, AgentID, Price)

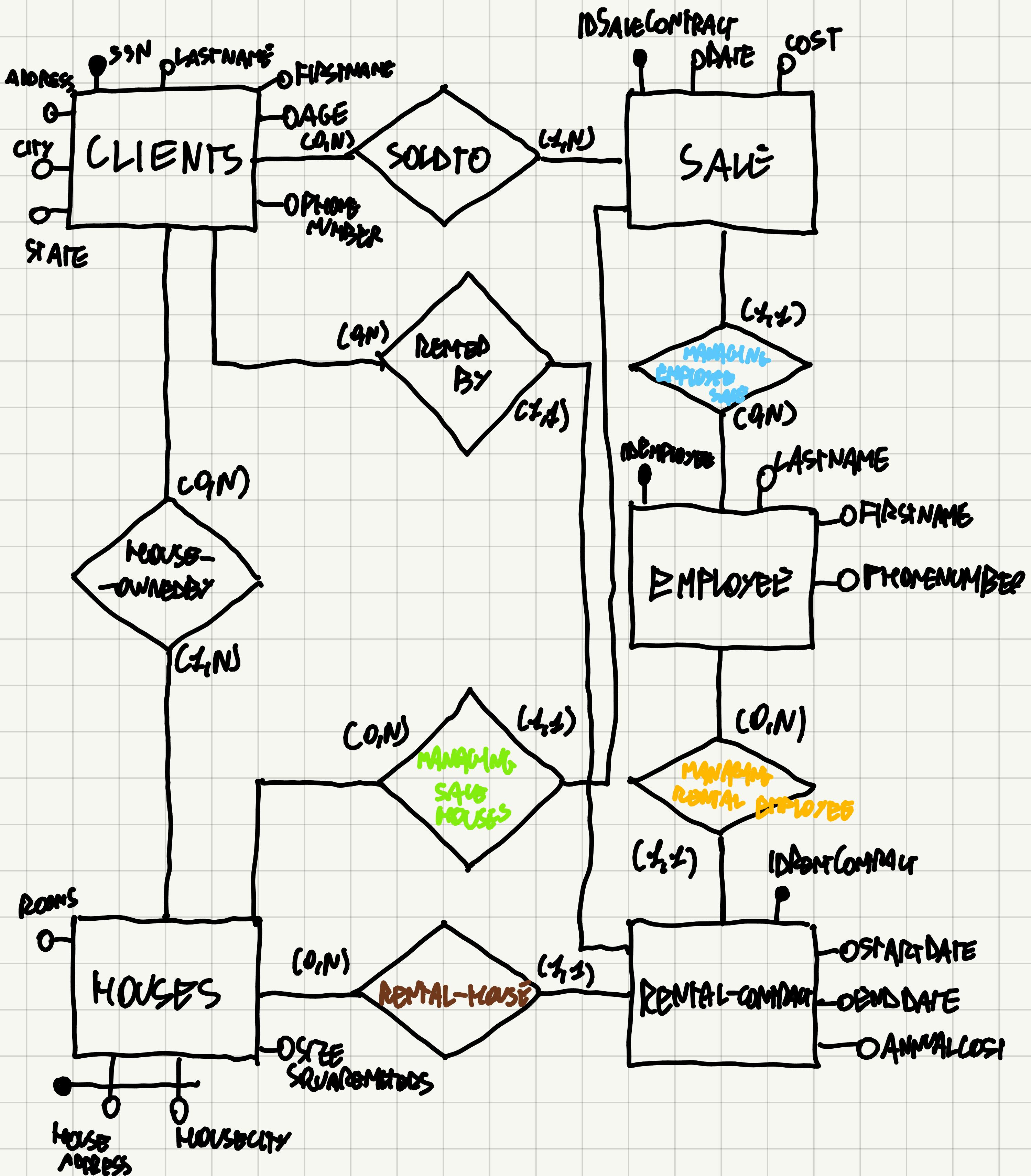
### **EXAMPLE OF WEAK ENTITY, DRAWN AS**

1. Provide, **for each** input data source, the reverse engineering from the logical to the conceptual schema (ER graph). (5 points)
2. Design an integrated global conceptual schema (ER graph) for *USARealEstateCompany* capturing **all** the data coming from both *LALuxuryHouses* and *USA Houses*, and provide the corresponding logical schema. (8 points)
3. Consider the query Q “Find the name and surname of the buyers who live in the city of Los Angeles and have bought at least one house larger than 100 square meters located in the city of Beverly Hills”.
  - a. Write GAV mappings between the schema of *USARealEstateCompany* and the two sources either in Datalog or SQL. **Write the mappings for the tables used to answer to query Q.** (3 points)
  - b. Consider query Q posed on *USARealEstateCompany*’s schema and write it either in Datalog or SQL. (3 points)
  - c. Show the rewriting of Q on the data sources either in Datalog or SQL. (3 points)

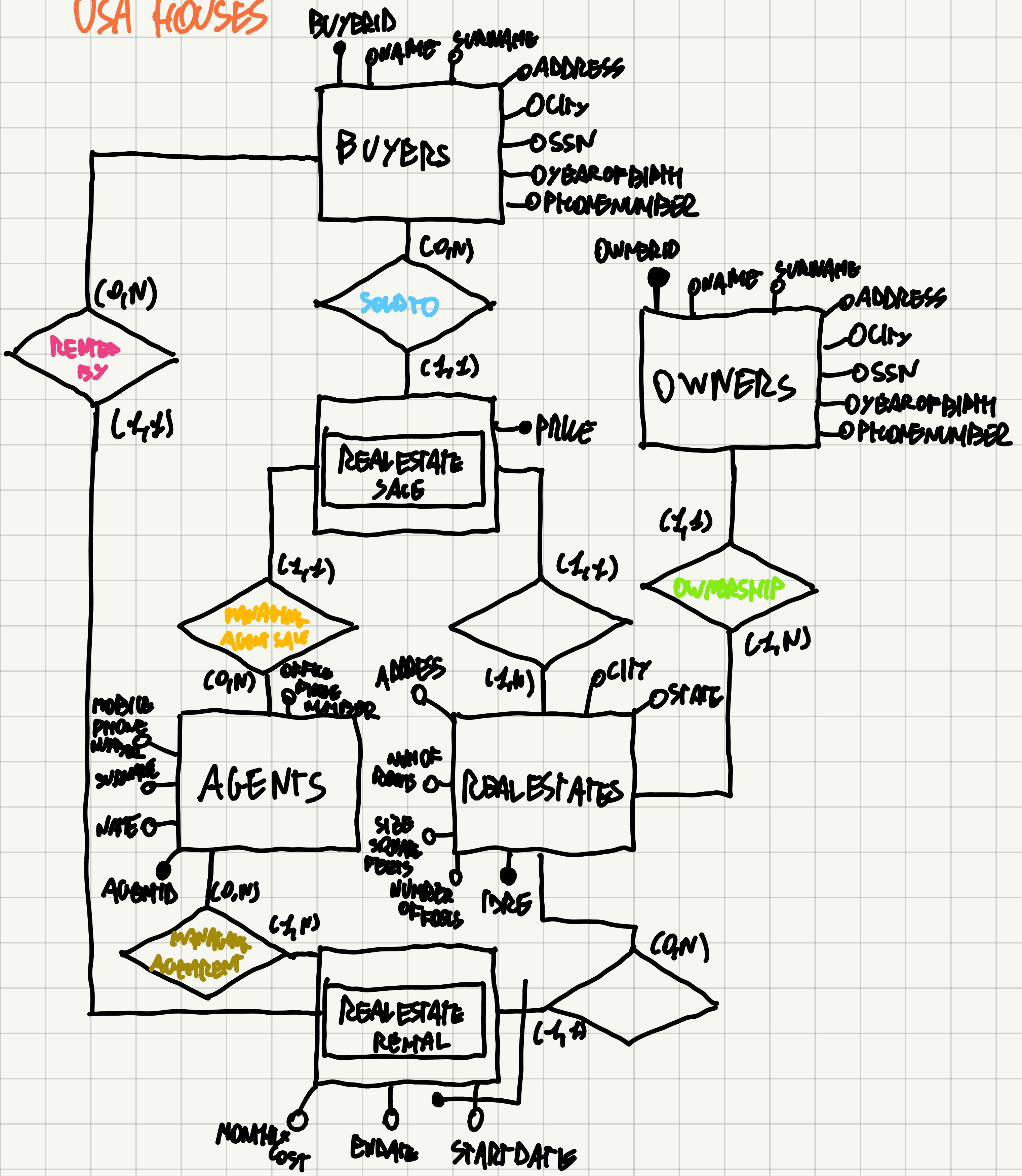
**Important:** 1) Spell out all your assumptions.

- 2) Avoid information loss as much as possible when defining the new schema.
- 3) List clearly all conflicts you detect during schema integration, if any.

# 1) LA LUXURY HOUSES



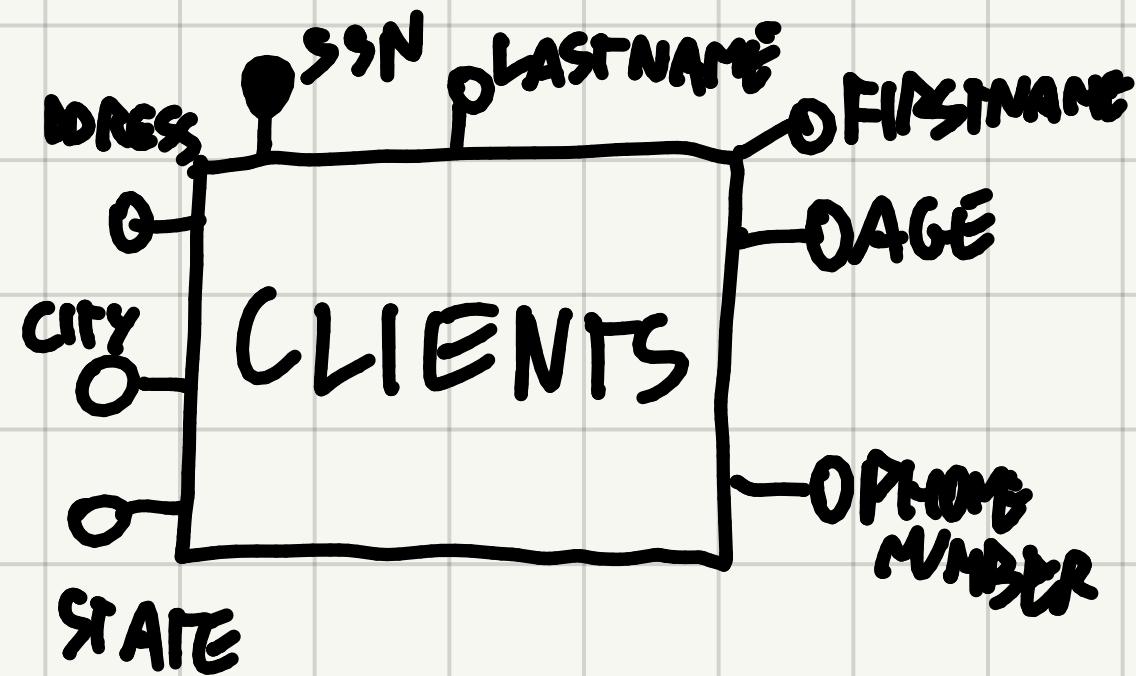
# USA HOUSES



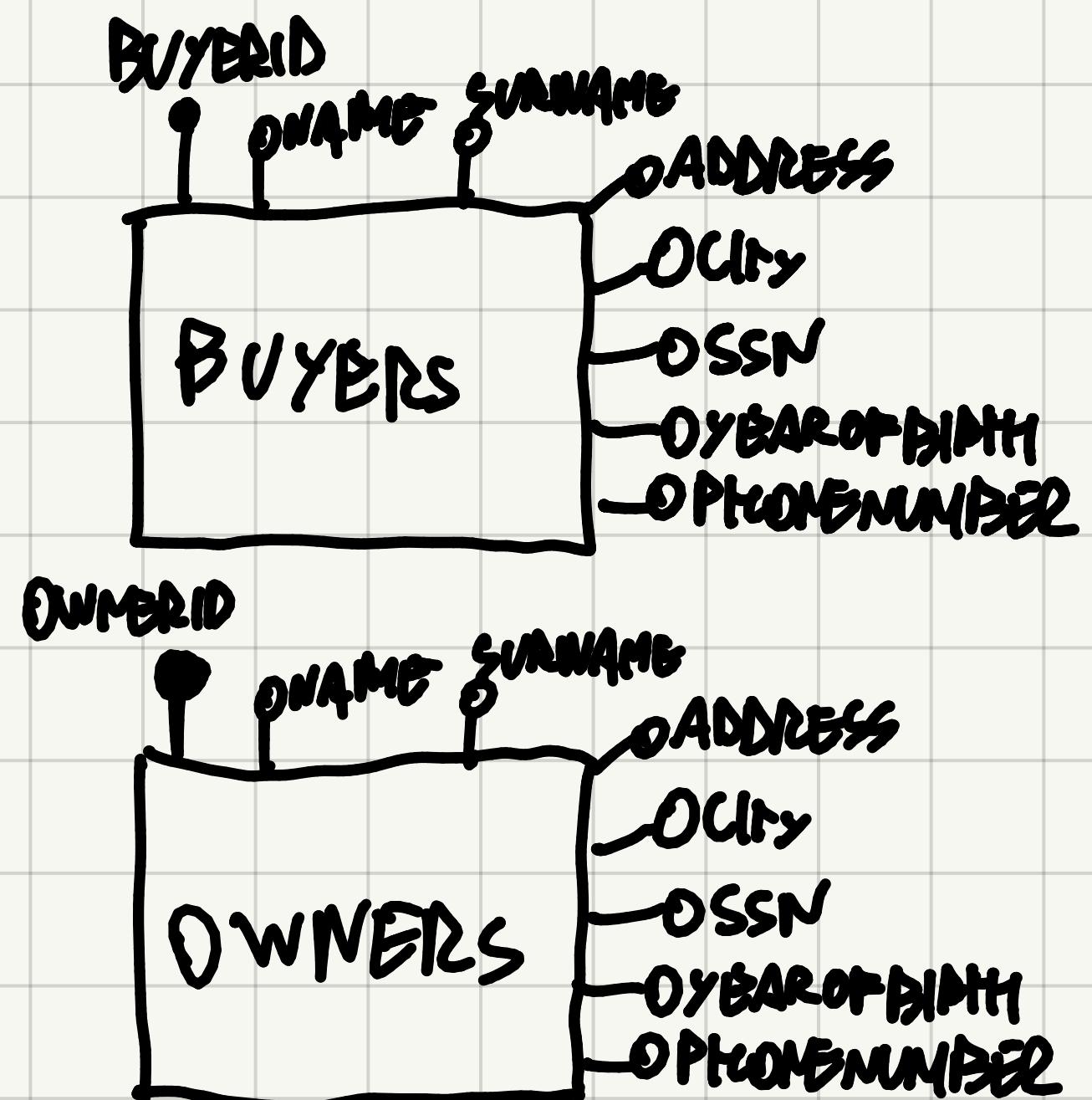
2)

## TABLES

LA Luxury Houses



USA Houses



BEST PRACTICE: CHOOSE THE MOST GENERAL ENTITY. IN THIS CASE, WE **UNITE**

**BUYERS AND OWNERS IN A GENERAL ENTITY **CLIENTS****  
**CONFLICTS**

**ENTITIES** (**CLIENTS; BUYERS, OWNERS**) → **CLIENTS**

**ATTRIBUTES** (**LAST NAME; SURNAME**) → **LAST NAME**  
**(FIRST NAME; NAME)** → **FIRST NAME**

\* THE CHOICE IS NOT IMPORTANT, EVEN "SURNAME - NAME" IS AN  
 ACCEPTABLE SOLUTION. THE IMPORTANT IS TO REMAIN COHERENT.

IN GENERAL, IT IS PREFERRED TO CHOOSE THE NAME OF THE ATTRIBUTE

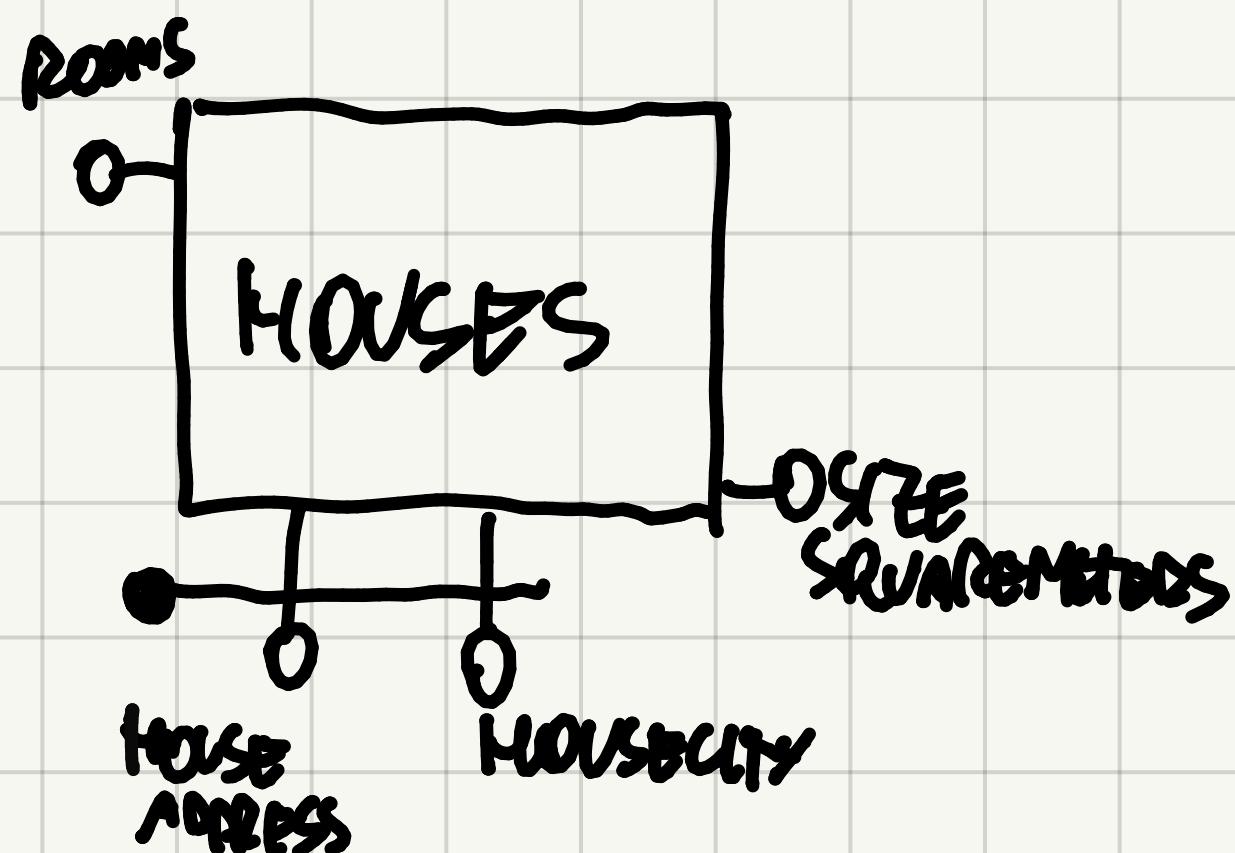
THAT MATCHES ITS RESPECTIVE ENTITY, AS WE DID NOW

**DATA SEMANTIC** (**AGE; YEAROFBIRTH**) → **YEAROFBIRTH**

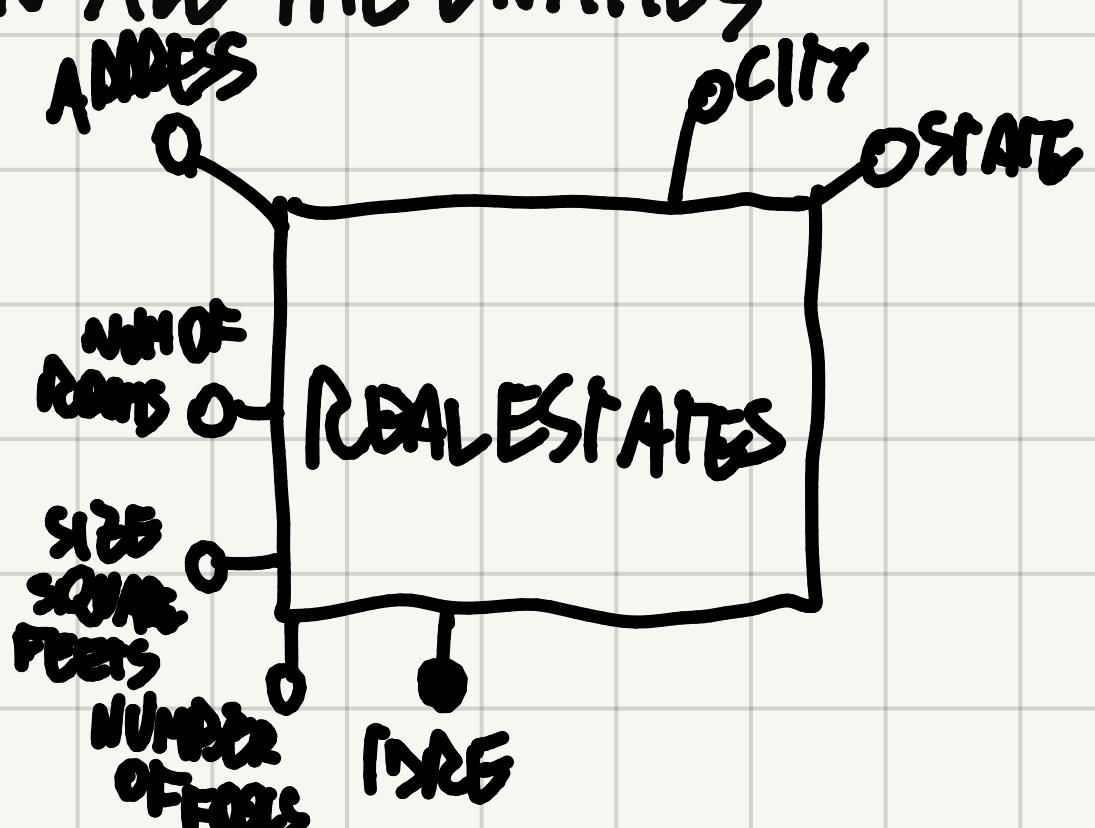
IT IS A BETTER CHOICE BECAUSE, WHILE A PERSON'S AGE CHANGES EVERY YEAR,  
 HIS YEAR OF BIRTH DO NOT AND, SO, IT DOESN'T NEED CONTINUOUS UPDATES

KEY CONFLICT (SSN; BUYERID, OWNERID) → SSN

WE NOTICE THAT SSN COMPARES IN ALL THE ENTITIES



CONFLICTS



SOLUTION

ENTITIES (HOUSES; REALESTATE) → REALESTATE

NAME

ATTRIBUTES (HOUSEADDRESS, HOUSECITY; ADDRESS, CITY) → ADDRESS, CITY

(ROOMS; NUMBEROFROOMS) → ROOMS

DATA SEMANTIC CONFLICTS (SIZE SQUAREMETERS; SIZE SQUARE FEET)

\*

→ SIZE SQUAREMETERS

\* AN ADVICE IS TO SEE IF THERE ARE SOME SPECIFIC REQUIREMENTS

FROM THE EXERCISE'S TRAK. IN THIS CASE, THE REQUESTED

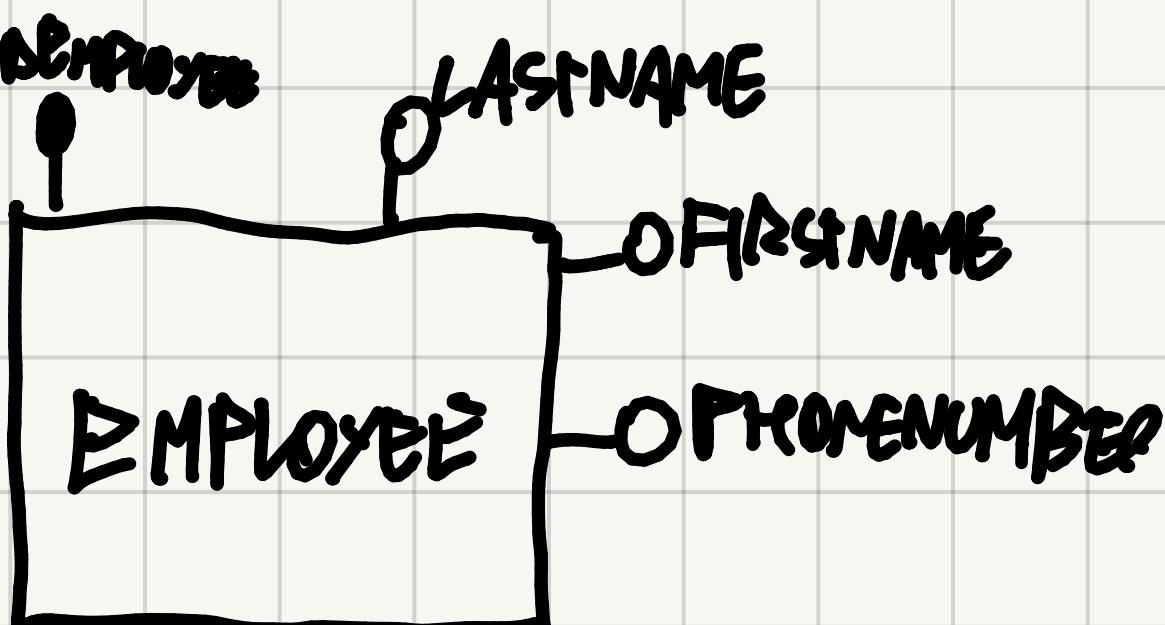
QUERY CALLS FOR "... SQUARE METERS..."

KEY CONFLICT (HOUSEADDRESS, HOUSECITY; IDRE) → ADDRESS, CITY

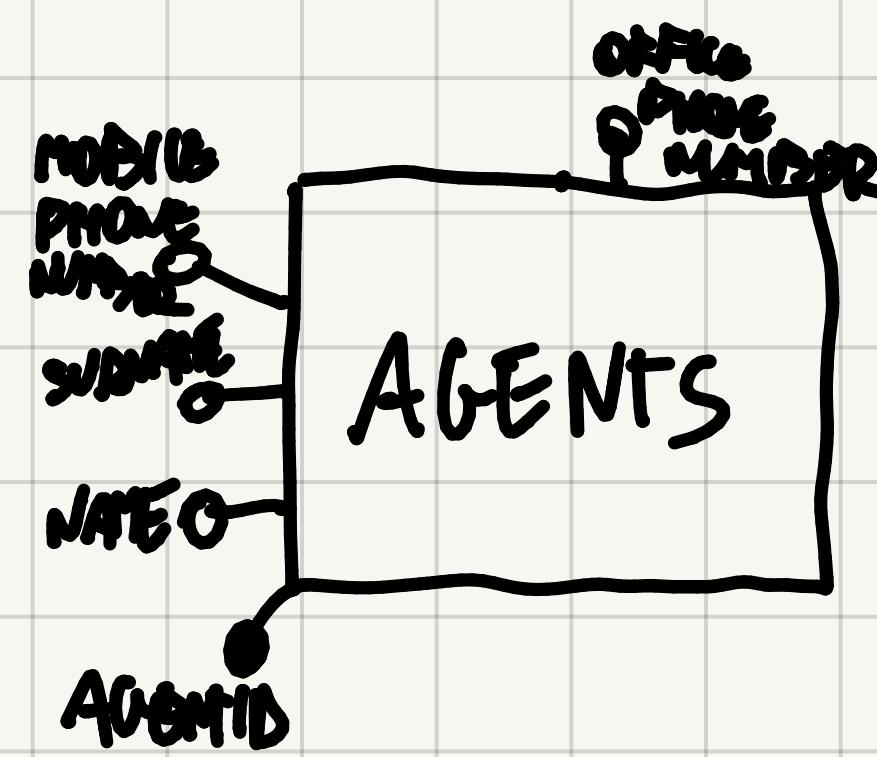
CARDINATALITY CONFLICTS

- NUMBER OF OWNERS (1..N); (1..1) → (1..N)

WE PREFER TO STAY AS MUCH GENERAL AS POSSIBLE.



CONFFLICTS



SOLUTION

ENTITIES (EMPLOYEE; AGENTS)



EMPLOYEE \*

\* IN THIS CASE, BOTH CHOICES ARE ACCEPTABLE

ATTRIBUTES (IDEMPLOYEE; AGENTID)



EMPLOYEE ID

(LASTNAME; SURNAME)



LASTNAME

(FIRSTNAME; NAME)



FIRSTNAME

(PHONE NUMBER; MOBILEPHONE NUMBER;  
OFFICEPHONE NUMBER)

OFFICEPHONE NUMBER

(WE ASSUME THAT PHONE NUMBER = OFFICEPHONE NUMBER)



CONFFLICTS



SOLUTION

NAME

ENTITIES (SALE, REAL ESTATE NAME)



SALE

ATTRIBUTES (COST; PRICE)



COST

KEY CONFLICT (IDSALECONTRACT, DATE + IDRE) →

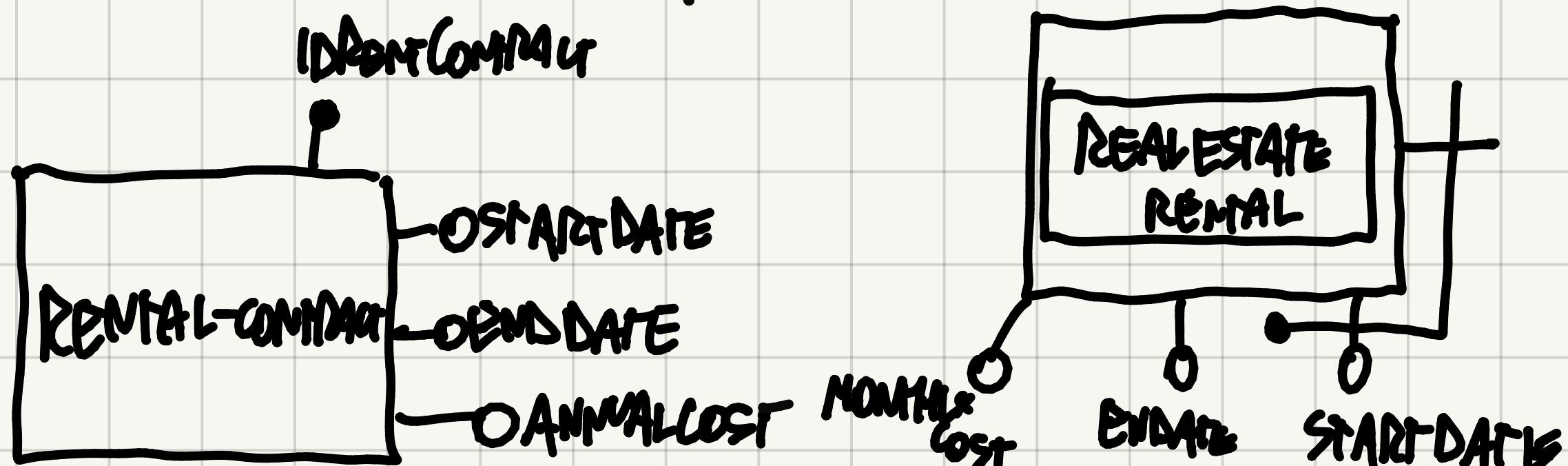
IDSALECONTRACT

[KEYGEN(ID, LA LUXURY HOUSES); KEYGEN(DATE+IDRE, USA HOUSES)] \*

\* WE CAN GENERATE A PRIMARY KEY FROM THE UNION OF THOSE

CARDINALITY CONFLICTS

- NUMBER OF BUYERS ( $(C_0, N)$ ;  $(1, \infty)$ )  $\longrightarrow (C_0, N)$



CONFFLICTS

NAME

SOLUTION

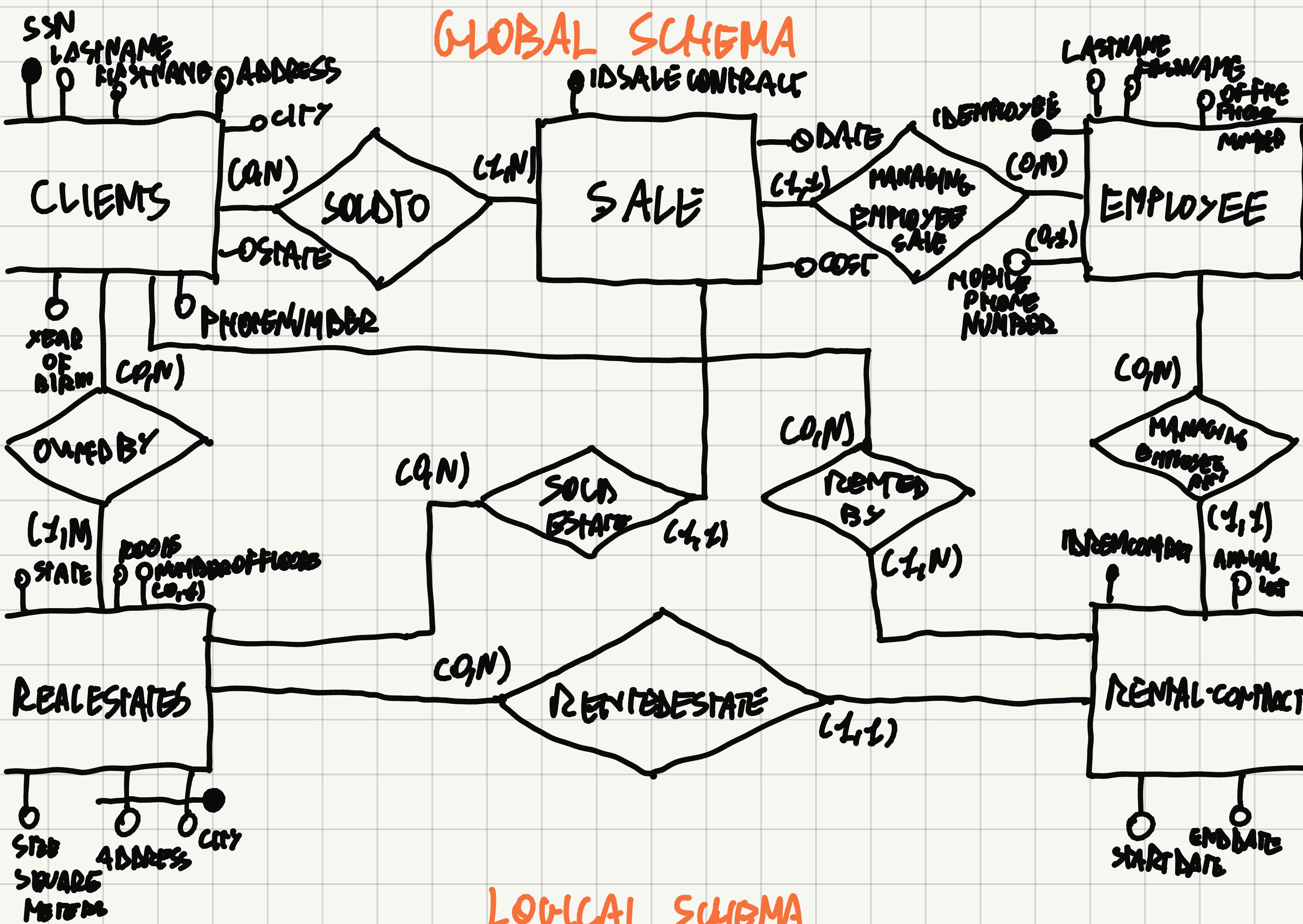
ENTITIES ( $RENTAL-CONTRACT$ ;  $REAL ESTATE REMAL$ )  $\longrightarrow RENTAL-CONTRACT$

DATASEMANTIC CONFLICTS (ANNUALCOST; MONTHLYCOST)  $\longrightarrow$  ANNUAL COST

KEY CONFLICT (IDRENTCONTRACT; IDRE + STARTDATE)  $\longrightarrow$  IDRENTCONTRACT

CARDINALITY CONFLICT

- NUMBER OF RENTERS ( $(C_0, N)$ ;  $(1, \infty)$ )  $\longrightarrow (C_0, N)$



**LOCAL SCHEMA**

**CLIENTS**(SSN, LASTNAME, FIRSTNAME, ADDRESS, CITY, STATE, YEAROFBIRTH,

**SALE**(IDSALECONTRACT, DATE, COST, ESTATEADDRESS, ESTATECITY, IDEMPLOYEE)

**OWNEDBY**(IDSALECONTRACT, CLIENT)

**REALESTATE**(ADDRESS, CITY, STATE, SIZESQUAREMETERS, ROOMS, NUMBEROFFLOORS\*)

**RENTAL-CONTRACT**(IDRENTCONTRACT, STARTDATE, ENDDATE, ANNUALCOST, IDEMPLOYEE,

ESTATEADDRESS, ESTATECITY)

**RENTEDBY**(IDRENTCONTRACT, CLIENT)

**OWNEDBY**(ESTATEADDRESS, ESTATECITY, CLIENT)

**MOBILEPHONENUMBER\***

**EMPLOYEE**(IDEMPLOYEE, LASTNAME, FIRSTNAME, OFFICEPHONENUMBER,

)

\* OPTIONAL ATTRIBUTE