

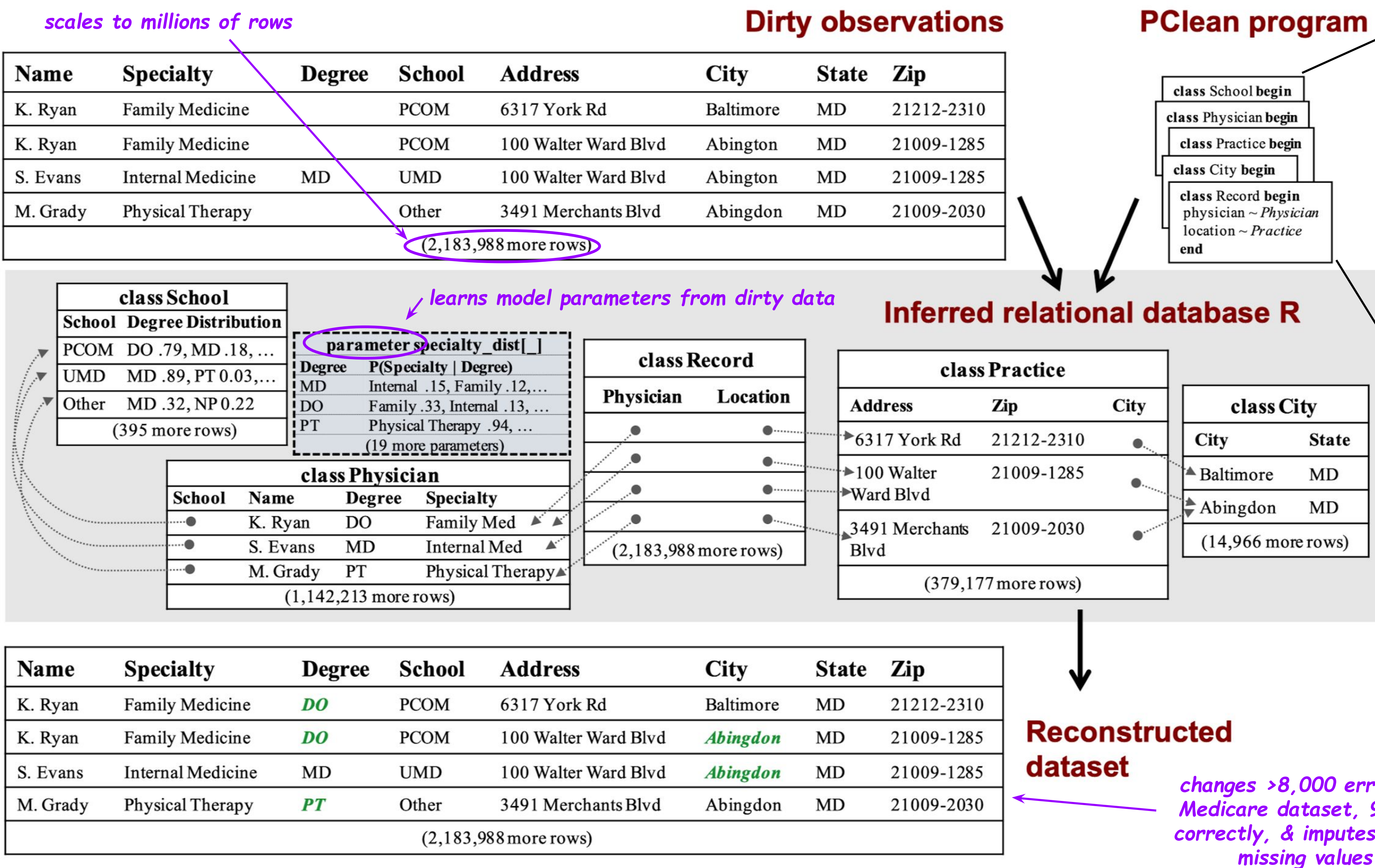


PClean: Bayesian Data Cleaning at Scale via Domain-Specific Probabilistic Programming

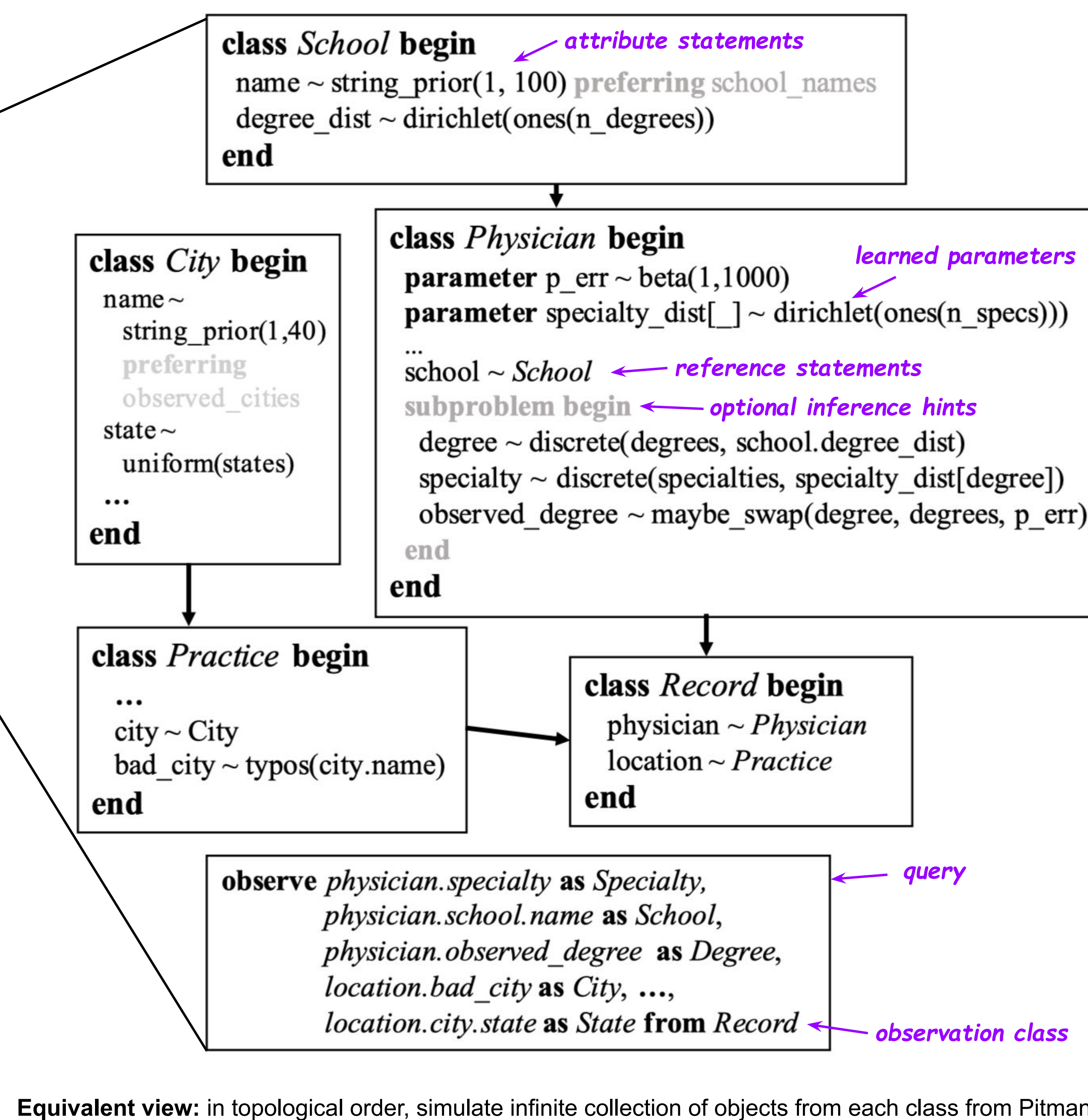
Alexander K. Lew, Monica Agrawal, David Sontag, Vikash K. Mansinghka
MIT

1. Overview

PClean is a domain-specific probabilistic programming language for inferring ground-truth relational databases from flat, dirty datasets.



2. PClean Modeling Language



Equivalent view: in topological order, simulate infinite collection of objects from each class from Pitman-Yor

Left: A PClean program defines a relational schema for a database of objects underlying the dirty data, along with a probabilistic relational model over object attributes.

Below: PClean uses a domain-general non-parametric structure prior over the number of objects of each class, and over their relationships.

GENERATESKELETON(\mathcal{C} , $|\mathbf{D}|$):

- ▷ Create one C_{obs} object per observed record $\mathbf{S}_{C_{obs}} := \{1, \dots, |\mathbf{D}|\}$
- ▷ Generate a class *after* all referring classes:

```

for class  $C' \in \text{TopoSort}(\mathcal{C} \setminus \{C_{obs}\})$  do
  ▷ Collect references to class  $C'$ 
   $\text{Refs}(C') := \{(r, Y) \mid r \in \mathbf{S}_{C'}, T(C'.Y) = C'\}$ 
  ▷ Generate targets of those references
   $\mathbf{S}_{C'} \sim \text{GENERATEOBJECTSET}(C', \text{Refs}(C'))$ 
  ▷ Assign reference slots pointing to  $C'$  for object  $r' \in \mathbf{S}_{C'}$  do
    for referring object  $(r, Y) \in r'$  do
       $r.Y := r'$ 
    end
  end
end
  
```

▷ Return the skeleton

return $\{\mathbf{S}_C\}_{C \in \mathcal{C}}, (r, Y) \mapsto r.Y$

GENERATEOBJECTSET(C , $\text{Refs}(C)$):

- $s_C \sim \text{Gamma}(1, 1)$; $d_C \sim \text{Beta}(1, 1)$
- ▷ Partition $\text{Refs}(C)$ into disjoint co-referring subsets; each represents an object $\mathbf{S}_C \sim \text{CRP}(\text{Refs}(C), s_C, d_C)$

3. PClean Inference Engine

1) Initializes latent database with *per-observation* SMC; fixes mistakes with *per-object* rejuvenation

- Non-parametric prior admits a **sequential representation** (right), enabling SMC that incorporates one observation at a time
- Exchangeability of CRPs enables ***per-object* rejuvenation moves**: choose any object, and propose new parameters and reference slots — possibly creating new objects of other classes as their targets

2) Fast, data-driven SMC / block rejuvenation proposals created just-in-time by a *proposal compiler*

- Translates subproblem to a Bayes net
- For each pattern of missingness in the data, compiles efficient enumeration-based SMC proposals
- As needed, compiles enumeration-based rejuvenation proposals
- When possible exploits conjugacy for continuous parameters

3) User-specified *inference hints* help scale to large datasets and variable domains

- subproblem begin ... end:** Group adjacent statements into a *subproblem* which becomes an intermediate target in SMC. Smaller subproblems lead to more scalable enumerative proposals, at some cost to proposal quality
- $\mathbf{x} \sim \mathbf{d}(\mathbf{E}, \dots, \mathbf{E})$ preferring \mathbf{E} :** Specify a dynamically computed list of values on which posterior is likely to concentrate; proposal enumerates these, and the enumerative proposal is mixed with prior using an adaptive mixture weight

GENERATEDATASET(Π , \mathbf{Q} , $|\mathbf{D}|$):

```

 $\mathbf{R}^{(0)} \leftarrow \emptyset$  ▷ Initialize empty database
for observation  $i \in \{1, \dots, |\mathbf{D}|\}$  do
   $\Delta_i^{\mathbf{R}} \leftarrow \text{GENERATEINCREMENT}(\mathbf{R}^{(i-1)}, C_{obs})$ 
   $\mathbf{R}^{(i)} \leftarrow \mathbf{R}^{(i-1)} \cup \Delta_i^{\mathbf{R}}$ 
   $r \leftarrow$  the unique object of class  $C_{obs}$  in  $\Delta_i^{\mathbf{R}}$ 
   $d_i \leftarrow \{X \mapsto r.Q(X), \forall X \in \mathcal{A}(\mathbf{D})\}$ 
end
return  $\mathbf{R} = \mathbf{R}^{(|\mathbf{D}|)}$ ,  $\mathbf{D} = (d_1, \dots, d_{|\mathbf{D}|})$ 
  
```

GENERATEDBINCR($\mathbf{R}^{(i-1)}$, root class C):

```

 $\Delta \leftarrow \emptyset$ ;  $r_* \leftarrow$  a new object of class  $C$ 
for each reference slot  $Y \in \mathcal{R}(C)$  do
   $C' \leftarrow T(C.Y)$ 
  for each object  $r \in \mathbf{R}_{C'}^{(i-1)} \cup \Delta_{\mathbf{R}_{C'}}$  do
     $n_r \leftarrow |\{r' \mid r' \in \mathbf{R}^{(i-1)} \cup \Delta \wedge \exists \tau, r'.\tau = r\}|$ 
     $r_*.Y \leftarrow r$  w.p.  $\propto n_r - d_{C'}$ , or  $\star$  w.p.  $\propto s_{C'} + d_{C'} |\mathbf{R}_{C'}^{(i-1)} \cup \Delta_{\mathbf{R}_{C'}}|$ 
    if  $r_*.Y = \star$  then
       $\Delta' \leftarrow \text{GENERATEDBINCR}(\mathbf{R}^{(i-1)} \cup \Delta, C')$ 
       $\Delta \leftarrow \Delta \cup \Delta'$ 
       $r_*.Y \leftarrow$  the unique  $r'$  of class  $C'$  in  $\Delta'$ 
    end
  end
for each  $X \in \mathcal{A}(C)$ , in topological order do
   $r_*.X \sim \phi_{C,X}(\cdot \mid \{r_*.U\}_{U \in Pa(C,X)})$ 
end
return  $\Delta \cup \{r_*\}$ 
  
```

Above: Unlike in general open-universe PPLs, PClean's model admits a **sequential representation**, in which one observation (and all *new* latent objects it is based on) is instantiated at each time step, enabling sequential Monte Carlo inference

4. Experiments & Results

PClean is expressive enough to model data cleaning benchmarks in <50 LoC, and unlike generic PPL inference (left), achieves SOTA accuracy + runtime compared to weighted logic + machine learning baselines.

Task	Metric	PClean	HoloClean (Unpublished)	HoloClean	NADEEF	NADEEF + Manual Java Heuristics
Flights	F_1	0.90	0.64	0.41	0.07	0.90
	Time	3.1s	45.4s	32.6s	9.1s	14.5s
Hospital	F_1	0.91	0.90	0.83	0.84	0.84
	Time	4.5s	1m 10s	1m 32s	27.6s	22.8s
Rents	F_1	0.69	0.48	0.48	0	0.51
	Time	1m 20s	20m 16s	13m 43s	13s	7.2s

Table 1: Results of PClean and various baseline systems on three diverse cleaning tasks.

