

Databáze ◊ poznámky k přednášce

## 7. Dokumentový model databáze

verze z 4. listopadu 2024

Začneme motivací. Představme si, že vlastníme jediný film *Dracula* z roku 1992 od režiséra *Francis Ford Coppola* a zajímáme se o herce *Gary Oldman*, *Anthony Hopkins* a *Tom Cruise*. Vydělujeme část reálného světa určenou vlastněnými filmy a herci, o které se zajímáme.

### 1 Schémata

Vhodné řetězce písmen anglické abecedy a podtržítka nazýváme **atributy**. Například `title`, `year`, `actor`, `name`.

Atributy rozdělíme na **atomické** a **kolekční**. Například můžeme určit `title`, `year` a `name` za atomické atributy a `actor` za kolekční atribut.

Každému atomickému atributu  $y$  je přiřazena spočetná množina  $D_y$  nazývaná **doména** atributu  $y$ . Například shodnou doménou  $D_{\text{title}}$  a  $D_{\text{name}}$  atributů `title` a `name` může být množina všech řetězců nad vhodnou abecedou a doménou  $D_{\text{year}}$  atributu `year` množina všech celých čísel.

Zavedeme pojem (**dokumentového**) **schématu hloubky**  $k$ , kde  $k$  je nezáporné celé číslo. Prázdná množina je schéma hloubky nula. Vezměme schémata  $S_1, \dots, S_n$  hloubky ostře menší než  $k$  a po dvou různé atributy  $y_1, \dots, y_n$  takové, že pokud je  $y_i$  atomický atribut, tak  $S_i$  je prázdná množina. Pak je množina

$$\{\langle y_1, S_1 \rangle, \dots, \langle y_n, S_n \rangle\}$$

schéma hloubky nejvýše  $k$ . Například  $\{\langle \text{title}, \emptyset \rangle, \langle \text{year}, \emptyset \rangle\}$  je schéma hloubky jedna a  $\{\langle \text{title}, \emptyset \rangle, \langle \text{actor}, \{\langle \text{name}, \emptyset \rangle\}\rangle\}$  je schéma hloubky dva.

Množina  $S$  je (**dokumentové**) **schéma**, jestliže existuje  $k$  takové, že  $S$  je schéma hloubky  $k$ . Například  $\emptyset$ ,  $\{\langle \text{title}, \emptyset \rangle, \langle \text{year}, \emptyset \rangle\}$  a  $\{\langle \text{title}, \emptyset \rangle, \langle \text{actor}, \{\langle \text{name}, \emptyset \rangle\}\rangle\}$  jsou schémata.

Dvojici  $\langle y, T \rangle$  ve schématu budeme také zapisovat jako „ $y.T$ “. Dále zápis „ $y.\emptyset$ “ budeme zkracovat na  $y$ . Schémata z předchozího příkladu můžeme tedy zkráceně zapsat  $\{\text{title}, \text{year}\}$  a  $\{\text{title}, \text{actor}.\{\text{name}\}\}$

Pro schéma  $S$  označíme  $Y(S)$  množinu

$$\{y \mid \langle y, S' \rangle \in S\}$$

atributů, které se nacházejí v schématu  $S$ . Například  $Y(\{\text{title}, \text{actor}.\{\text{name}\}\}) = \{\text{title}, \text{actor}\}$ .

Všimněme si, že schéma  $S$  je zobrazení přiřazující atributu z  $Y(S)$  schéma  $S(y)$  nižší hloubky. Například pro  $S = \{\text{title}, \text{actor}.\{\text{name}\}\}$  je  $S(\text{actor}) = \{\text{name}\}$ .

## 2 Kolekce

Jisté dvojice  $\langle y, d \rangle$ , kde  $y$  je atribut, nazýváme **komponenty**. Definujeme je pro atomické a kolekční atributy zvlášť.

Jestliže je  $y$  atomický atribut a  $d \in D_y$ , pak  $\langle y, d \rangle$  je (**atomická**) **komponenta**. Například  $\langle \text{title}, \text{"Dracula"} \rangle$  nebo  $\langle \text{year}, 1992 \rangle$  jsou atomické komponenty.

Nechť  $S$  je schéma. Pak množinu komponent

$$t = \{\langle y_1, t_1 \rangle, \dots, \langle y_n, t_n \rangle\}$$

nežveme **dokumentem** nad  $S$ , jestliže

1.  $y_1, \dots, y_n$  jsou po dvou různé,
2.  $\{y_1, \dots, y_n\} = Y(S)$ ,
3. pro každý kolekční atribut  $y \in Y(S)$  je  $t(y)$  kolekce nad  $S(y)$ .

Například

$$\{\langle \text{title}, \text{"Dracula"} \rangle, \langle \text{year}, 1992 \rangle\}$$

je dokumentem nad  $\{\text{title}, \text{year}\}$ .

Konečnou množinu dokumentů nad  $S$  nazýváme **kolekcí** nad  $S$ . Například:

$$C_2 = \{\{\langle \text{name}, \text{"Gary Oldman"} \rangle\}, \{\langle \text{name}, \text{"Anthony Hopkins"} \rangle\}\}$$

je kolekce nad  $\{\text{name}\}$ . Prázdnou množinu nazýváme též **prázdnou kolekcí**.

Jestliže  $y$  je kolekční atribut a  $C$  je kolekce nad  $S$ , pak  $\langle y, C \rangle$  je (**kolekční**) **komponenta**. Například:

$$\langle \text{actor}, \{\{\langle \text{name}, \text{"Gary Oldman"} \rangle\}, \{\langle \text{name}, \text{"Anthony Hopkins"} \rangle\}\}\rangle$$

je kolekční komponenta.

Dostáváme, že množina

$$\begin{aligned} &\{\langle \text{title}, \text{"Dracula"} \rangle, \\ &\langle \text{actor}, \{\{\langle \text{name}, \text{"Gary Oldman"} \rangle\}, \\ &\quad \{\langle \text{name}, \text{"Anthony Hopkins"} \rangle\}\}\rangle \end{aligned}$$

je dokumentem nad  $\{\text{title}, \text{actor}.\{\text{name}\}\}$ .

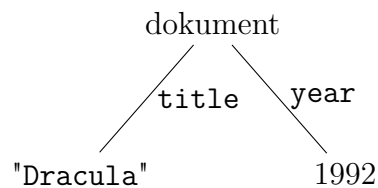
Všimněme si, že dokument nad  $S$  je zobrazení, které atributu  $y$  z  $Y(S)$  přiřazuje **hodnotu atributu**  $t(y)$ . Atributům  $Y(S)$  říkáme **atributy dokumentu**. Například výše uvedený dokument má atribut `title` s hodnotou "Dracula".

Dokument  $t$  můžeme přirozeně zobrazit stromovým diagramem tak, že kořen zobrazuje dokument  $t$ . Následníci uzlu zobrazujícího dokument jsou v jednoznačné korespondenci s jeho komponentami. Hrana mezi uzlem a následníkem s přiřazenou komponentou je popsána atributem komponenty. Uzel s přiřazenou atomic-kou komponentou zobrazuje její hodnotu. Uzel s přiřazenou kolekční komponentou zobrazuje její kolekci. Následníci uzlu zobrazujícího kolekci jednoznačně odpovídají prvkům kolekce. Následník uzlu zobrazujícího kolekci zobrazuje odpovídající prvek kolekce.

Například dokument

$$\{\langle \text{title}, \text{"Dracula"} \rangle, \langle \text{year}, 1992 \rangle\}$$

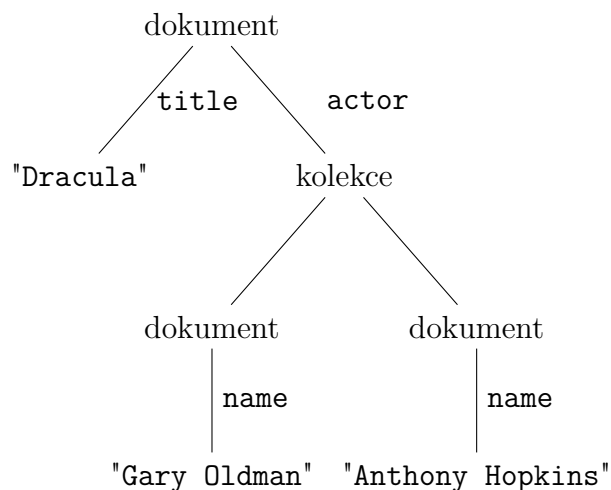
zobrazíme diagramem:



a

$$\{\langle \text{title}, \text{"Dracula"} \rangle, \langle \text{actor}, \{\{\langle \text{name}, \text{"Gary Oldman"} \rangle\}, \{\langle \text{name}, \text{"Anthony Hopkins"} \rangle\}\}\rangle\}$$

diagramem:



K textovému zápisu dokumentů budeme používat formát JSON (JavaScript Object Notation). Atomické hodnoty zapíšeme přímo. Například řetězec "Dracula" nebo číslo 1992. Dokument  $\{\langle y_1, d_1 \rangle, \dots, \langle y_n, d_n \rangle\}$  zapíšeme řetězcem:

```
{
  "y1": d1,
  ⋮
  "yn": dn
}
```

Například dokument

```
{⟨title, "Dracula"⟩, ⟨year, 1992⟩}
```

zapišeme řetězcem:

```
{
  "title": "Dracula",
  "year": 1992
}
```

Kolekci  $\{t_1, \dots, t_n\}$  zapišeme řetězcem:

```
[
  t1,
  ⋮
  tn
]
```

Například kolekci:

```
{{⟨name, "Gary Oldman"⟩}, {⟨name, "Anthony Hopkins"⟩}}
```

zapišeme řetězcem:

```
[
  {
    "name": "Gary Oldman"
  },
  {
    "name": "Anthony Hopkins"
  }
]
```

a konečně dokument:

```
{⟨title, "Dracula"⟩,
  ⟨actor, {{⟨name, "Gary Oldman"⟩},
           {⟨name, "Anthony Hopkins"⟩}}⟩}
```

zapišeme řetězcem:

```

{
  "title": "Dracula",
  "actor": [
    {
      "name": "Gary Oldman"
    },
    {
      "name": "Anthony Hopkins"
    }
  ]
}

```

Množinu všech dokumentů nad  $S$  označíme  $\text{Doc}(S)$ .

### 3 Charakteristické vlastnosti

Výroková forma  $V(t)$ , kde  $t$  je proměnná nabývající hodnot z množiny  $\text{Doc}(S)$  se nazývá **vlastnost** nad  $S$ . Například „ $t(\text{name})$  je jméno herce“ je vlastnost nad  $\{\text{name}\}$ . Proměnnou  $t$  většinou můžeme vynechat a zkráceně psát například jen „ $\text{name}$  je jméno herce“.

Vlastnost  $V$  nad  $S$  je **charakteristická vlastnost kolekce**  $C$  nad  $S$ , jestliže  $t \in C$ , právě když  $V(t)$  pro každý dokument  $t$  nad  $S$ . Zkráceně budeme říkat, že  $C$  **je určeno**  $V$ .

Platí, že  $V$  je charakteristická vlastnost  $C$ , právě když

$$C = \{t \in \text{Doc}(S) \mid V(t)\}.$$

Například kolekce

```

[
  {
    "title": "Dracula",
    "actor": [
      {
        "name": "Gary Oldman"
      },
      {
        "name": "Anthony Hopkins"
      }
    ]
  }
]

```

má charakteristickou vlastnost  $V(t)$ :

Řetězec  $t(\mathbf{title})$  je jméno filmu a  $t(\mathbf{actor})$  je kolekce s charakteristickou vlastností  $V(t')$  rovnou „Herec  $t'(\mathbf{name})$  hrál ve filmu  $t(\mathbf{title})$ “.

Aby bylo možné zjednodušit vlastnost. Zavedeme operaci přejmenování.

**Operace přejmenování.** Atomické atributy  $y_1$  a  $y_2$  jsou **stejného typu**, jestliže  $D_{y_1} = D_{y_2}$ .

Vezměme dvě schémata  $S_1$  a  $S_2$  a bijekci  $h: Y(S_1) \rightarrow Y(S_2)$ , která splňuje, že

1.  $S_1(y) = S_2(h(y))$  pro každý  $y \in Y(S_1)$ ,
2.  $y$  a  $h(y)$  jsou stejného typu pro každý atomický atribut  $y \in Y(S_1)$ .

Pak bijekce  $h$  je **přejmenování atributů** mezi  $S_1$  a  $S_2$ . Například pro  $S_1 = \{\mathbf{name}\}$  a  $S_2 = \{\mathbf{move.name}\}$  je  $h = \{\langle \mathbf{name}, \mathbf{movie.name} \rangle\}$  přejmenování atributů.

Nechť  $h$  je přejmenování atributů mezi  $S_1$  a  $S_2$  a  $t$  dokument nad  $S_1$ .

Pak **přejmenování dokumentu**  $t$  podle  $h$  je dokument

$$\rho_h(t) = \{\langle h(y), d \rangle \mid \langle y, d \rangle \in t\}$$

nad  $S_2$ . Například pro výše uvedené  $h$  a dokument  $t = \{\langle \mathbf{name}, \text{"Gary Oldman"} \rangle\}$  je  $\rho_h(t) = \{\langle \mathbf{movies.name}, \text{"Gary Oldman"} \rangle\}$ .

**Přejmenování kolekce**  $C$  nad  $S_1$  podle  $h$  je kolekce

$$\rho_h(C) = \{\rho_h(t) \mid t \in C\}$$

nad  $S_2$ .

Předpokládejme, že  $Y(S_1) = \{y_1, \dots, y_n, x_1, \dots, x_m\}$ , kde  $h(x_i) = x_i$  pro každé  $1 \leq i \leq m$ . Pak  $\rho_h$  můžeme značit  $\rho_{h(y_1) \leftarrow y_1, \dots, h(y_n) \leftarrow y_n}$ . Přejmenování podle výše uvedeného  $h$  můžeme značit  $\rho_{\mathbf{actors.name} \leftarrow \mathbf{name}}$ .

**Příklad vlastnosti.** Vrátime se k příkladu charakteristické vlastnosti. Přejmenujeme atributy kolekce  $t(\mathbf{actor})$ :

Řetězec  $t(\mathbf{title})$  je jméno filmu a  $\rho_{\mathbf{actor.name} \leftarrow \mathbf{name}}(t(\mathbf{actors}))$  je kolekce s charakteristickou vlastností  $V(t')$  rovnou „Herec  $t'(\mathbf{actor.name})$  hrál ve filmu  $t(\mathbf{title})$ “.

Nyní můžeme proměnné  $t, t'$  odstranit:

Řetězec  $\mathbf{title}$  je jméno filmu a  $\mathbf{actor}$  je kolekce s charakteristickou vlastností „Herec  $\mathbf{actor.name}$  hrál ve filmu  $\mathbf{title}$ “.

Vlastnost zjednodušíme:

Řetězec  $\mathbf{title}$  je jméno filmu a  $\mathbf{actor}$  je určeno vlastností: herec  $\mathbf{actor.name}$  hrál ve filmu  $\mathbf{title}$ .

Běžně budeme vlastnost kolekce a v ní vložené vlastnosti uvádět zvlášť. Tedy charakteristická vlastnost kolekce je:

Řetězec `title` je jméno filmu.

a charakteristická vlastnost kolekčního atributu `actor` je:

Herec `actor.name` hrál ve filmu `title`.

Obecně při formulaci charakteristické vlastnosti nad  $S$  s proměnnou  $t$  provedeme implicitně přejmenování  $\rho_h(t(y))$  pro každý kolekční atribut  $y \in Y(S)$ , kde  $h$  přejmenuje atribut  $y' \in S(y)$  na  $y.y'$ .

## 4 Kolekční proměnné

Pokud má dokument atribut `_id`, říkáme, že **má identifikátor**. Hodnota atributu `_id` se nazývá **identifikátor dokumentu**.

Uvažujme schéma  $S$  takové, že `_id`  $\in Y(S)$ . **Kolekční proměnná** nad schématem  $S$  je proměnná, jejíž hodnotou je kolekce nad  $S$ . Dokumenty v  $S$  musí mít po dvou různé identifikátory. Schéma  $S$  nazýváme **typem proměnné**. Můžeme mít například kolekční proměnnou `movie` nad `{_id, title, actor.{name}}`.

Každá kolekční proměnná má určenou **charakteristickou vlastnost**. Vždy platí, že hodnota proměnné má její charakteristickou vlastnost. Například `movie` má charakteristickou vlastnost:

„Film `_id` se jmenuje `title`“

a `actor` je určeno vlastností:

„Herec `actor.name` hrál ve filmu `title`.“

Hodnota proměnné `movie` je kolekce:

```
[
  {
    "_id": 1,
    "title": "Dracula",
    "actor": [
      {
        "name": "Gary Oldman"
      },
      {
        "name": "Anthony Hopkins"
      }
    ]
  }
]
```

Budeme pracovat s databázovým systémem MongoDB. Návod na jeho instalaci naleznete v dokumentu `07_mongodb.pdf`. Proměnné ani jejich schémata nemusíme v systému deklarovat. Každá kolekční proměnná *collection* má na začátku hodnotu rovnou prázdné kolekci.

Příkaz:

```
db.collection.insertMany(C)
```

změní hodnotu kolekční proměnné *collection* na  $C \cup C'$ , kde  $C'$  je původní hodnota proměnné *collection*. Musí platit, že množiny  $C$  a  $C'$  jsou disjunktní. Například:

```
db.movie.insertMany([
  {
    "_id": 1,
    "title": "Dracula",
    "actor": [
      {
        "name": "Gary Oldman"
      },
      {
        "name": "Anthony Hopkins"
      }
    ]
  }
])
```

**Kolekční výraz** je výraz, jehož hodnota je kolekce. Hodnota kolekčního výrazu:

```
db.collection.find()
```

je rovna hodnotě proměnné *collection*. Například:

```
> db.movie.find()
[
  {
    _id: 1,
    title: 'Dracula',
    actor: [ { name: 'Gary Oldman' }, { name: 'Anthony Hopkins' } ]
  }
]
```

Systém reprezentaci JSON zjednodušuje tak, že atributy není potřeba uzavírat do dvojitéch uvozovek a řetězce mohou být uzavřeny mezi jednoduché uvozovky.

Můžeme přidat další dokumenty do kolekce:

```
db.movie.insertMany([
  {
    _id: 2,
    title: 'Oppenheimer',
    actor: [
```



```

    {
      name: 'Gary Oldman'
    }
  ]
}, {
  _id: 3,
  title: 'The Matrix',
  actor: []
}])

```

Nová hodnota proměnné:

```
> db.movie.find()
```

```

[
  {
    _id: 1,
    title: 'Dracula',
    actor: [ { name: 'Gary Oldman' }, { name: 'Anthony Hopkins' } ]
  },
  { _id: 2, title: 'Oppenheimer', actor: [ { name: 'Gary Oldman' } ] },
  { _id: 3, title: 'The Matrix', actor: [] }
]

```

Příkaz:

```
db.collection.drop()
```

nastaví proměnnou *collection* na prázdnou kolekci. Například:

```
> db.movie.drop()
```

## Otázky a úkoly na cvičení

1. Vymezme si část reálného světa určenou režiséry, jejichž filmy vlastníme. U režiséra nás zajímá rok jeho narození a jeho filmy, které vlastníme. U filmu pak chceme znát jeho název a rok vydání. Navrhněte schéma  $S$  pro uložení popsaných dat.
2. Vezměte schéma  $S$  z prvního úkolu a určete množinovým zápisem kolekci  $C$  popisující následující stav vymezené reality. Od režiséra Tim Burton narozeného v roce 1958 vlastníme film *Edward Scissorhands* z roku 1990 a film *Alice in Wonderland* z roku 2010 a od režiséra Francis Ford Coppola narozeného v roce 1939 vlastníme film *Megalopolis* z roku 2024.
3. Vezměme kolekci  $C$  z druhého úkolu a zapište ji ve formátu JSON.

4. Deklarujte kolekční proměnnou nad schématem  $S$  doplněným o atribut `_id` z prvního úkolu a určete její charakteristickou vlastnost.
5. Změňte hodnotu proměnné ze čtvrtého úkolu tak, aby odpovídala stavu reality popsaném v druhém úkolu.