

# Obalovač (Wrapper) web stránok

---

- Čo sú a na čo slúžia obalovače?
- Problémy
  - ◆ parsovanie HTML dokumentov
  - ◆ verifikácia korektnosti obalovača
  - ◆ čitateľnosť a rozšíriteľnosť obalovača
  - ◆ navigácia
- Existujúce prostriedky

# Prečo potrebujeme obalovače - Aplikácie

---

- Porovnávanie ponúk/cien
- Vyhľadávanie
- Adaptácia starých alebo nevhodných rozhraní
- Integrácia údajov (napr. údaje z pobočiek)
- Sledovanie konkurencie

# Vstup obalovačov

The screenshot displays the Quanta web editor interface. The title bar shows the file path: `file:///home/gyuri/Desktop/program.html - Quanta`. The menu bar includes: File, Edit, View, Bookmarks, Project, Toolbars, DTD, Tags, Plugins, Tools, Window, Settings, Help. The toolbar contains various icons for file operations and editing. On the left, the Document Structure pane shows a tree view of the HTML document:

- Resources (link) [XHTML 1.0 Strict]
- Images [XHTML 1.0 Strict]
- Links (anchor) [XHTML 1.0 Strict]
- Document Structure
  - DTD block
    - DOCTYPE html PUBLIC "-//W3C//DT
  - html
    - head
      - title: Personalized Web
      - meta
      - meta
      - link
      - link
    - body
      - div
        - div
          - h1: PeWe (Personalized Web)
          - ul
            - href http://www.
            - href http://www.
            - href http://www.
        - div
        - ul
        - div
        - div
        - div
        - P

The main editing area shows a preview of the rendered HTML page. The content includes:

- A main heading: **PeWe (Personalized Web) Group**
- A sub-heading: [UI SI @ FIIT @ STU Bratislava](#)
- A bulleted list of links:
  - ◆ [Introduction](#)
  - ◆ [Members](#)
  - ◆ **Program**
  - ◆ [Current projects](#)
  - ◆ [Results](#)

Below the preview, the source code editor shows the following HTML code:

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN"
"http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">

<html xmlns="http://www.w3.org/1999/xhtml" xml:lang="sk" lang="sk">

<head>
<title>Personalized Web</title>
<meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
<meta name="language" content="sk" />
<link rel="stylesheet" type="text/css" media="screen, projection" href="master.css" />
<link rel="shortcut icon" href="img/favicon.ico" />
</head>

<body><div id="container">
<div id="caption">
<h1>PeWe (Personalized Web) Group</h1>
<p><a title="Institute of Informatics and Software Engineering"
href="http://www.uisi.fiit.stuba.sk">UISI</a> @ <a title="Faculty of Informatics and Information
```

The status bar at the bottom indicates: Ready. | INS | Line: 1 Col: 1

# Prístupy wrappovania

---

- Práca s HTML zdrojovým kódom ako so stringom  
(regulárne výrazy)
- Využitie stromovej štruktúry HTML kódu
- Vytváranie wrappera vymenovaním príkladov
- Adaptívne prístupy

# HTML dokumenty

---

- HTML dokumenty sú pološtrukturované dáta, kde stromová štruktúra tagov môže prezradiť niečo o význame dát
- Často však HTML dokumenty nespĺňajú špecifikáciu
  - ◆ využitie robustných parserov (mozilla, beautifulsoap a iné)
  - ◆ “očistenie” HTML dokumentov, ich konverzia do XHTML

# Vývoj a udržiavanie obalovača

---

- Vývojový nástroj má zadať jednoduchý rámec na tvorbu obalovača
- Nástroj má poskytovať spôsob na verifikáciu: ak obalovač zlyhá kvôli zmene stránky, vývojár má o tom dostať správu
- Identifikácia prípadov použitia, model používateľa  
– aké znalosti sa od používateľa očakávajú

# Navigácia

---

- Aby sme sa dostali k údajom, potrebujeme sa prebiť cez linky, cookies, heslá, formuláre
- Je potrebné riešiť autentifikáciu na webové sídla
- Zachytenie a nakódenie navigačných sekvencií môže viesť neprehľadnému kódu
- Spôsoby na zachytenie navigácie:
  - ◆ Proxy
  - ◆ Zmena liniek v HTML kóde
  - ◆ Modifikovaný prehliadač (resp. plug-in)

# Existujúce prostriedky a služby

---

## ■ I-point

- ◆ ponúkajú konkrétne wrapovanie konkrétnych sídiel
- ◆ neponúkajú nástroj na tvorbu

## ■ Kapowtech

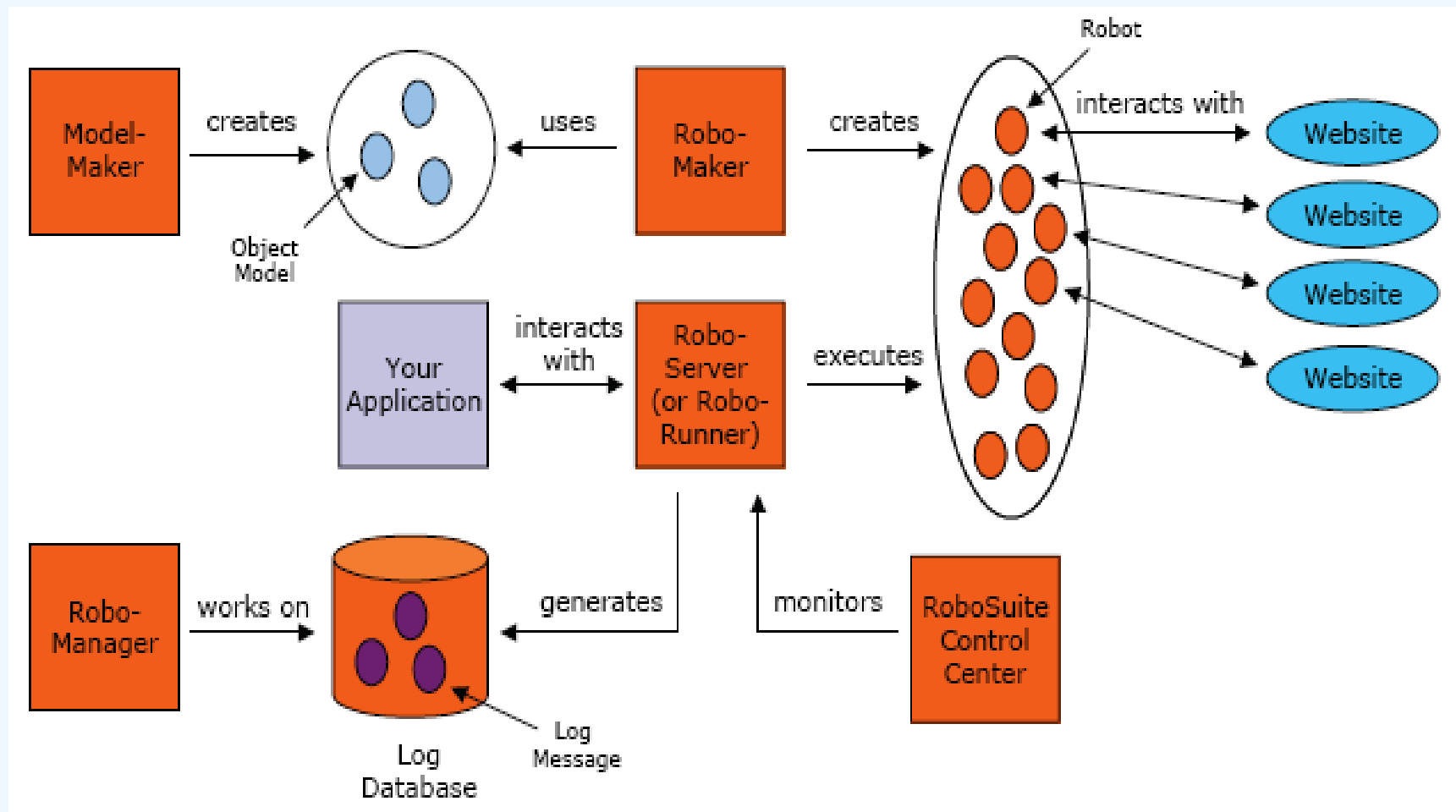
- ◆ integrovaný nástroj umožňujúci jednoduchú tvorbu obalovačov
- ◆ Wrapper je automat generujúci XML

## ■ LiXto

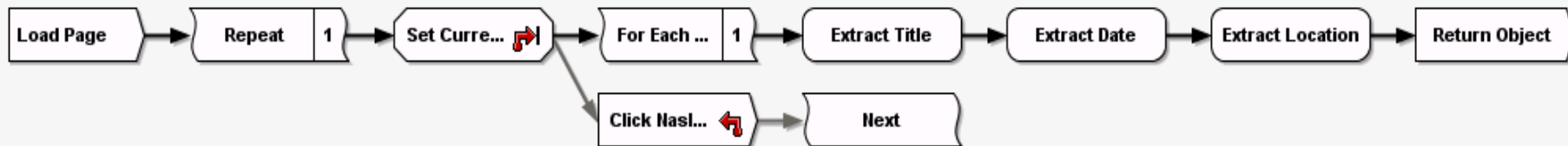
- ◆ dva nástroje: Visual Wrapper a Transformation Server
- ◆ Wrapper je sada pravidiel – rule-based prístup



# Kapowtech



# Kapowtech



Počet aktuálnych ponúk: 5290

1. **Obchodník**

24.10.2005, LOGOLOGIC SLOVAKIA, spol. s r.o.  
Bratislavský kraj, Bratislava, Bratislava-Hattalova 12,  
BA a okolí

2. **Finance Manager english**

24.10.2005, HR Partners, s.r.o.

99. **Administratívny pracovník, logistika, marketing**

24.10.2005, EMM International, spol. s r.o.  
Bratislava, Dlhé diely

100. **Hosteska/Promotérka**

24.10.2005, iLeo/DMMS, spol. s r.o.  
Slovenská republika

Strana: <<prvá | <predchádzajúca | 2 | nasledujúca >> | posledná >>

Nenašli ste vyhovujúcu pracovnú ponuku? Chcete, aby zamestnávateľia našli vás? Zvýšte šancu nájdenia práce a vložte do našej databázy svoj životopis.

- Visual Wrapper

- ◆ Vizuálny nástroj na tvorbu obalovača

- Transformation Server

- ◆ Nástroj na transformáciu, prezentáciu a doručenie výsledkov získaných z obalovačov

# Otázky

---

- Aký je vhodný model používateľa obalovača?
- Aké sú bežné a najčastejšie problémy pri tvorbe obalovača?
- Aký je vhodný jazyk na zápis obalovača?
- ...