



M Ű E G Y E T E M 1 7 8 2

BUDAPESTI MŰSZAKI ÉS GAZDASÁGTUDOMÁNYI EGYETEM

VILLAMOSMÉRNÖKI ÉS INFORMATIKAI KAR
IRÁNYÍTÁSTECHNIKA ÉS INFORMATIKA TANSZÉK

A WIKIPÉDIA VALÓSIDEJŰ NYELVI FELDOLGOZÁSA OSGI ALAPOKON

SZAKDOLGOZAT

Készítette UNICSOVICS MILÁN GYÖRGY *Konzulens* HÉDER MIHÁLY, SIMON BALÁZS

2013. december 7.

Tartalomjegyzék

Kivonat	4
Abstract	5
Bevezető	6
1. Hasonló megoldások vizsgálata	7
1.1. Egy ígéretes megoldás: Wikipedia Miner	8
1.2. Tanulságok	9
2. Az OSGi keretrendszeréről	11
2.1. Bundle	11
2.2. Service, Service Registry	13
2.3. Életciklus	15
3. Tervezés	16
3.1. WikiBot bundle	17
3.2. Parser bundle	18
3.2.1. Wikipédia dump feldolgozás	21
3.3. Databaseconnector bundle	21
3.4. Database	22
3.5. Logger bundle	23
3.6. Statistics bundle	24
4. Implementáció	28
4.1. WikiBot bundle	28
4.2. Parser bundle	29
4.3. DatabaseConnector bundle	32
4.4. Logger bundle	33
4.5. Statistics bundle	34
5. Tesztelés és mérések	35
5.1. A teszteléshez használt eszközök	35
5.2. A feldolgozólánra épülő kutatómodul	36
5.3. Mérési eredmények	37

6. Értékelés, továbbfejlesztési lehetőségek	41
Ábrák jegyzéke	44
Táblázatok jegyzéke	45
Forráskódok jegyzéke	46
Irodalomjegyzék	47
Függelék	49
6.1. WikiProcessor	49

HALLGATÓI NYILATKOZAT

Alulírott *Unicsovics Milán György*, szigorló hallgató kijelentem, hogy ezt a szakdolgozatot meg nem engedett segítség nélkül, saját magam készítettem, csak a megadott forrásokat (szakirodalom, eszközök stb.) használtam fel. Minden olyan részt, melyet szó szerint, vagy azonos értelemben, de átfogalmazva más forrásból átvettem, egyértelműen, a forrás megadásával megjelöltem.

Hozzájárulok, hogy a jelen munkám alapadatait (szerző(k), cím, angol és magyar nyelvű tartalmi kivonat, készítés éve, konzulens(ek) neve) a BME VIK nyilvánosan hozzáférhető elektronikus formában, a munka teljes szövegét pedig az egyetem belső hálózatán keresztül (vagy autentikált felhasználók számára) közzétegye. Kijelentem, hogy a benyújtott munka és annak elektronikus verziója megegyezik. Dékáni engedéllyel titkosított diplomatervek esetén a dolgozat szövege csak 3 év eltelte után válik hozzáférhetővé.

Budapest, 2013. december 7.

Unicsovics Milán György
hallgató

Kivonat

A szemi-struktúrált adatok felhasználása rendkívül népszerű a mesterséges intelligencia, természetes nyelvi feldolgozás, de az informatika más területein is. Az ilyen szemi-struktúrált adatforrásokat, mint például a legnagyobb kollaboratívan szerkeszthető tudásbázist, a Wikipédiát felhasználó megoldások száma mégis igen kevés. Szakdolgozatom célja egy kutatók számára is használható rendszer tervezésének és fejlesztésének bemutatása, mely a Wikipédiát dolgozza fel, és teszi elérhetővé kutatások céljából. A rendszer egy feldolgozólánc formájában fog megvalósulni, mely az OSGi keretrendszert használva, annak flexibilitását kihasználva egyszerűen kiegészíthető, illetve továbbfejleszthető lesz, így alkalmassá válik további kutatások kiindulópontjaként is szolgálni.

A tervezés előtt mindenképp meg kell vizsgálni a hasonló megoldásokat, hiszen azok tapasztalatait érdemes felhasználni: hibáikból tanulni lehet, a jól bevált megoldásokat pedig meg kell fontolni, hogy beépíthetők-e az alkalmazás architektúrájába. Az elérhető szemi-struktúrált adatforrásokat feldolgozó alkalmazásokat végigvizsgálva mindegyikben egy alapvető hiányosság tűnt fel: egyik sem elég flexibilis, és a legtöbb program csak egy nagyon speciális célra alkalmas, így általános kutatások alapjául nem használhatóak. A megoldások közös jellemzője, hogy alapvetően három jól meghatározott, nagyobb előfeldolgozó komponenst tartalmaznak, így ezeket a komponenseket (információk gyűjtése, feldolgozása, eltárolása) terveztem meg és implementáltam én is a követelményeknek megfelelő módon.

A rendszer architektúráját nagyban meghatározza a flexibilitást és karbantarthatóságot eredményező OSGi technológia, melynek komponens modellje egy megszokottól eltérő Java nyelvű fejlesztést tesz lehetővé.

A tényleges tervezésnél az OSGi komponenseket külön-külön készítettem el, figyelve arra, hogy az alkalmazás teljesítménye minél nagyobb legyen. A gyorsaságot legnagyobb mértékben a többszálú futásnak köszönheti az alkalmazás, de a feldolgozás több részének aszinkron kivitelezése és nagy teljesítményű technológiák alkalmazása is sokat segített.

Az implementációs fázisban ismertetem a kipróbált technológiákat, melyeknek milyen előnyei, hátrányai vannak, és miért esett rájuk a választás. A technológiai elemeken kívül a tervezési minták használata is hangsúlyos az elkészített munkámban.

A kivitelezés eredményeit végül egy mérésen keresztül mutatom be, mely fázis átvezet a továbbfejlesztési lehetőségekhez. Itt mérésekkel alátámasztva ismertetem, melyik részt, hogyan lehetne továbbfejleszteni a specifikációtól eltérve, hogy a követelményeknek megfelelően, de mégis egy hiánypótló alkalmazás szülessen a szemi-struktúrált feldolgozórendszerek körében.

Abstract

Using of semi-structured datas is very popular in the field of artificial intelligence (AI), natural language processing (NLP) but other branch of computer sciences too. Nevertheless, there are only a few solution that are using semi-structured datasources, for example the greatest collaboratively editable knowledge base, the Wikipedia. My thesis aims to present the full lifecycle of a new research tool, that can help in computer science researches and it is able to process the articles of Wikipedia and publishes the result for the researchers.

Before the design, similar solutions have to be investigated, because from their experience we can learn a lot. We can avoid common mistakes, but for example use a good architecture. After examining the available solutions it can be see, that none of them is flexible and they can be used only at specific tasks. Common property of the applications that they all have three major, well defined components (gathering, processing and storing the datas). These three components were designed and implemented in my thesis regarding the specified requirements.

A system's architecture based on OSGi module system and service platform that comes with flexibility and maintainability. The OSGi framework extends Java language's ability with a component driven development method.

At the design phase the separate OSGi components were created one-by-one taking care of high performance. The speed of the application comes from the multi threaded approach, but the async implementation and the usage of high performance technologies bring a lot of speed too.

At the implementation chapter the tested technologies are demonstrated, covering the advantages and disadvantages of each technology. Besides the technical elements, the design patterns are very important in my thesis too.

The results of the implementation will be demonstrated through a measurement, which chapter leads to opportunities of further development. The further development ideas are proved with measurements and it is explained how is it possible to modify the application to create a useful application in the field of semi-structured data processing systems.

Bevezető

Mai világunkban a tudás a legnagyobb érték. A tudás valójában kontextusba ágyazott információ, mely elemi adatokból épül fel. Hosszú ideje irányulnak kutatások és fejlesztések az informatikában a tudás minél hatékonyabb megszerzésére, azonban, ha egy nehéz és összetett problémát akarunk megoldani valamilyen módszerrel, azt tapasztaljuk, hogy a hozzá szükséges tudás megszerzése lesz mindig a szűk keresztmetszet (ez az ún. *knowledge acquisition bottleneck*).

Ez a helyzet egy ideje a különféle adatábrázolási módszereknek köszönhetően kezdett megváltozni. A struktúráltan nehezen feldolgozható adatok mellett megjelentek a struktúrált és szemi-struktúrált adatok, melyekből az információt sokkal gyorsabban és könnyebben lehet kinyerni, és olyan automatizált folyamatokban is felhasználhatóak már, ahol korábban emberi közreműködés lett volna szükséges. Megjelentek a kollaboratívan szerkeszthető tudásbázisok, ezek közül is a legnépszerűbb és legnagyobb a Wikipédia, mely egy szemi-struktúrált információforrás. Ezen tudásbázisokat főleg a mesterséges intelligencia [7], természetes nyelvi feldolgozás [3], valamint a számítógépes nyelvészet területén, de az informatika szinte minden ágában ugyanúgy felhasználgák.

A szemi-struktúrált információforrások, mint a Wikipédia (továbbiak például a Flickr, Twitter vagy a Yahoo! Answers) egyesítik a struktúrált és struktúráltatlan források előnyeit, így rendkívül jó kiindulópontjai különféle kutatásoknak. A struktúráltatlan adatokkal szembeni előnye, hogy gépek számára is értelmezhető információt tárolnak; a struktúrált adatok előállításánál és kezelésénél pedig kevésbé erőforrás-igényes a szemi-struktúrált adatok létrehozása és karbantartása.

Az informatika ezen területén történő kutatások már 2005–2006 körül megkezdődtek, azonban az ezek eredményeit felhasználó fejlesztések száma nem túl sok. Céloom tehát az eddig elkészült és elérhető fejlesztések vizsgálata, információk és tapasztalatok összesítése, valamint ezeket továbbfejlesztve egy korszerűbb rendszer összeállítása, melyre alapozva eredményes kutatásokat lehet kezdeményezni, a legnagyobb elérhető szemi-struktúrált erőforrás a Wikipédia segítségével, a lehető legflexibilisebb módon.

1. fejezet

Hasonló megoldások vizsgálata

Mielőtt a saját Wikipédia alapú rendszer tervezésének nekiláttam volna megvizsgáltam a hasonló célú, elérhető megoldásokat, hogy utána a tapasztalatokból kiindulva láthassak neki a munkának. Megpróbáltam összegyűjteni az összes hasonló témakört feldolgozó alkalmazást, majd ezeket egyesével megvizsgáltam, milyen előnyökkel, hátrányokkal rendelkeznek.

WikiNet[19] A WikiNet egy teljesen struktúrált adatszerkezetű tudásbázist, másnéven ontológiát épít fel. Ez egy Perl nyelvű szkriptek gyűjteményéből álló megoldás, a Wikipédia egy statikus, letöltött verzióját használja forrásként, a kinyert fogalmakat és kapcsolatokat külön szöveges fájlokban tárolja el.

Elérhetőség:

<http://www.h-its.org/english/research/nlp/download/wikinet.php>

DBpedia Ez a megoldás a Wikipédia dump-ján alapul, belőle tényszerű információkat nyer ki és rögzít struktúrált formában, melyet ezután a weben közzétéve a felhasználók számára egy SQL szerű nyelvvel lekérdezhetővé tesz. A alkalmazás Scala, Java nyelven íródott, adatbázisként Virtuoso Universal Server-t használ.

Elérhetőség:

<http://dbpedia.org/>

BabelNet[16] A BabelNet egy Java nyelven írt alkalmazás, ontológiát készít más adatforrások (Wikipédia dump, WordNet) alapján és ezekből egy „enciklopédikus szótárt” készít a felhasználók számára. Egy adott fogalomra keresve a szemantikailag kapcsolódó fogalmak is megjelennek, valamint ezek szinonímái (a BabelNet elnevezése szerint *synset*-ek), és minden szinonímához tartozik egy rövid definíció (*gloss*) több nyelven.

Elérhetőség:

<http://lcl.uniroma1.it/babelnet/>

Java Wikipedia Library[18] A JWPL egy Java alapú alkalmazás, mely egy interfészt ajánl ki, amivel a Wikipédia tartalmához lehet hozzáférni. A JWPL tartalmaz egy Mediawiki Markup parser-t, mellyel a letöltött Wikipédia dump-ot beolvassák, a wikitext-ből átalakított

szövegekből optimalizált adatbázisokat készítenek, melyekhez végül hozzáférési felületet nyújtanak.

Elérhetőség:

<http://www.ukp.tu-darmstadt.de/software/jwpl/>

Wikipedia Preprocessor Ezzel az eszközzel más alkalmazások számára lehet egy Wikipédia dump-ot feldolgozhatóbb formába hozni. A statikus dump-ból kiinduló alkalmazásoknál erre szükség is van, mivel a például Wikipédia 2007. július 19-i XML formátumban letölthető tartalma is 12 GB méretű, amit nyers formában szinte lehetetlen hatékonyan felhasználni.

Elérhetőség:

<http://www.cs.technion.ac.il/~gabr/resources/code/wikiprep/>

YAGO2[11] A YAGO2 a Wikipédia egy korábban letöltött tartalmán, és egyéb online tudásbázisokon (WordNet, GeoNames) alapuló ontológia. Az ontológiához több formátumban lehet hozzáférni, ezt kisebb Java nyelven írt konvertáló eszközökkel lehet megtenni.

Elérhetőség:

<http://www.mpi-inf.mpg.de/yago-naga/yago/>

1.1. Egy ígéretes megoldás: Wikipedia Miner

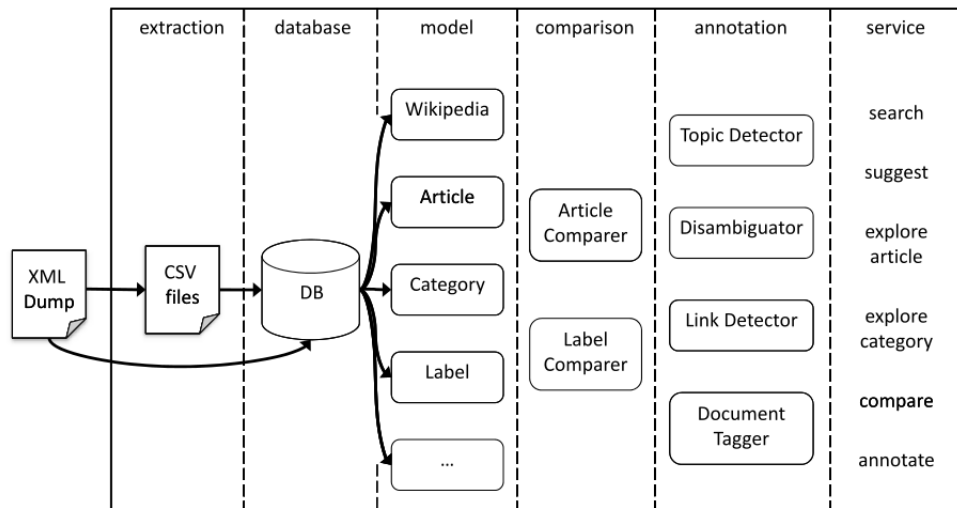
A Wikipedia Miner egy olyan nyílt forráskódú teljes egészében Java nyelven írt eszközkészlet, mellyel lehetővé válik kutatók és fejlesztők számára, hogy alkalmazásukban egyszerűen hozzáférjenek a Wikipédia tartalmához. A hozzáférési felületet Java API-n keresztül nyújtja az alkalmazás, a tudásbázis tartalmát pedig egy összegzett formában használja fel. Ezenfelül a rendszer olyan technológiákat használ fel, mint az elosztott számítási keretrendszert nyújtó Apache Hadoop és a Weka adatbányászati szoftver.

A Wikipedia Miner alkotói szerint három lehetőség van a szemi-struktúrált adatforrások felhasználására: vagy kész ontológiákat használunk, mint például a YAGO2, vagy egy nyers letöltött dump-ból elkészítjük a saját ontológiánkat, mint például a JWPL és az alapján kezdünk dolgozni; a harmadik lehetőség, hogy folyamatosan frissítjük az adatforrásunkat, így mindig naprakész forrást használunk.

A rendszer áttekintő architektúráját mutatja be a 1.1. ábra. A rendszer belépési pontjára érkezik egy nagy méretű XML formátumú fájl, mely a Wikipédia hivatalos API-ján keresztül lett letöltve és tartalmazza a naprakész Wikipédia teljes állományát.

Az *extraction* package funkciója, hogy kinyerjen az XML forrásból egy összegzett tartalmat. Az adatkinyerési folyamat használja a nagy méretű XML fájl feldolgozásához a Hadoop és MapReduce technológiákat, valamint a Google GFS fájlrendszerét. Az adatok kinyerése tehát elosztott rendszeren történik, a 27 GB méretű teljes angol Wikipédia feldolgozása a mérések szerint egy 2 magos, 2,66 GHz órajelű processzorral és 4 GB memóriával rendelkező 30 gépből álló clusternek 2,5 órájába telik.

A Wikipédia kivonatolása után a teljes XML fájl és az összegzett tartalmak is bekerülnek a *database* package adatbázisába. Adatbázisként Berkeley DB Java Edition-t használnak, mellyel



1.1. ábra. Wikipedia Miner architektúra (forrás: Artificial Intelligence, Wikipedia and Semi-Structured Resources[7])

akár egy teljes adatbázist is a memóriában lehet tartani, ezzel nagyon gyors lekérdezhetőséget elérve.

A Wikipedia Miner keretrendszert használók a tartalmakhoz a *model* package absztrakciós felületen keresztül férhetnek hozzá, mely becsomagolja a kinyert információkat és jól dokumentált osztályok formájában teszi közzé (mint például: Wikipedia, Article, Category).

A további csomagok (*comparison* és *annotation*) már a felhasználók számára használható szemantikus tartalom kezelését segítő eszközök, a *service* package-ben pedig konkrét szolgáltatások találhatók, melyeket felhasználhatnak a fejlesztők, illetve kiegészíthetik őket.

1.2. Tanulságok

Amint látható a fentebb leírt Wikipédiát felhasználó eszközök szinte mind kizárólag a Wikipédia egy korábban letöltött változatán (dump) alapulnak, mely módszernek több hátránya is van. A letöltött tudásbázis mérete rendkívül nagy, így azt kezelni nehézkes, feldolgozási ideje rendkívül hosszú (a Wikipedia Processor feldolgozási ideje például körülbelül 43,5 óra). A hosszú feldolgozási idő nem mindig engedhető meg, ráadásul amíg nincs feldolgozott forrás, a rendszer sem működőképes. Újabb dump letöltésekor kezdhető előlről a feldolgozás, így többször is megakaszthatja ez a folyamat a rendszer működését, a rendelkezésre állási időt csökkentve.

A másik megoldás, melyet a Wikipedia Miner (1.1. alfejezet) is demonstrál a Wikipédia folyamatos *on-the-fly* feldolgozása. Míg az előző módszer figyelmen hagyja a közösségi tudásbázisok fő erejét, hogy rendkívül dinamikusban fejlődnek, az *on-the-fly* megoldás kihasználja azt, és mindig a friss adatokkal dolgozik.

További fontos megállapítások, hogy egyik rendszer sem elég flexibilis: futásidőben új komponens beépítése egyáltalán nem lehetséges, nem lehet a feldolgozólánchoz új elemet (kutatást végző modult) illeszteni, szerkezetük szinte mindegyiknek rendkívül statikus és csak arra a célra használhatóak konkrétan, amilyen speciális feladat ellátására kitalálták. Ebből adódóan rendkívül

specifikusak és bonyolultak, így újabb kutatások indítása eredményeiket felhasználva nehézkes.

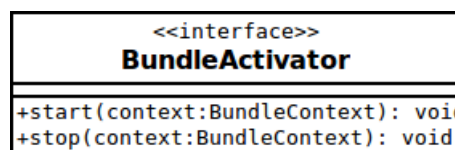
A szoftverek bármely módosítása újrafordítást, rendszerleállást eredményez, ami egy aktívan használt program számára nagy kiesést jelenthet. On-the-fly feldolgozásnál a rendszer leállása kihagyott, nem feldolgozott információkkal jár, így a tudásbázis töredezetté válhat. Ezen okok miatt szükséges egy olyan megoldás, mellyel a rendszer sokkal flexibilisebbé tehető és a fenti problémák áthidalhatóvá válnak.

2. fejezet

Az OSGi keretrendszeréről

Az OSGi Alliance [1] által fejlesztett OSGi (*Open Services Gateway initiative*) egy Java nyelvű, dinamikus, szolgáltatás orientált komponens modell és keretrendszer, mellyel kiterjeszthető az alap Java nyelven készült programok funkcionalitása. Használatával elérhető válik egy komponens alapú fejlesztés, ahol alkalmazásokat (illetve komponenseket) távolról elérve telepíthetjük, elindíthatjuk, leállíthatjuk, frissíthetjük, vagy törölhetjük anélkül, hogy a teljes alkalmazást leállítanánk és újra elindítanánk. A standard Java nyelven íródott alkalmazásoknál ez egy fontos hiányosság, hiszen az információs rendszerek sok területén nem engedhetők meg akár a pillanatnyi leállások sem.

A flexibilitás és újrafelhasználhatóság követelményeknek tehát nagyon jól megfelel az OSGi keretrendszer, ezért is esett rá a választás. Léteznek természetesen más, komponens alapú modellt használó technológiák, például: Microsoft Common Object Model, Enterprise JavaBeans, CORBA Component Model. Ezek közül azonban a legtöbb egy komplex programozási modellen alapul, amely megköveteli konvenciók követését és az OSGi-al szemben nem támogatják a komponensek dinamikus frissítését, így flexibilitás szempontjából ilyen módon elmaradnak.



2.1. ábra. BundleActivator osztálydiagram

2.1. Bundle

Az OSGi technológiát [2][14][17] használó alkalmazások kisebb komponensekre (az OSGi terminológiát használva csomagokra, azaz *bundle*-ökre) vannak bontva. Ezen csomagok elkészíthetők, lefordíthatók, telepíthetők egymástól függetlenül, életciklusukat maga az OSGi keretrendszer felügyeli. A *bundle*-ök gyakorlatilag nem mások, mint a jól ismert JAR fájlok (Java osztályok, és egyéb erőforrások becsomagolva), azzal a különbséggel, hogy a leíró manifest állományban a szabvány által kiterjesztett módon további fejlécek találhatók meg. Példát látunk manifest állományra a 2.1. kódrészletben. A manifest-ben lévő metaadatok nagy része embe-

ri felhasználás céljából szerepel és nem módosítják az OSGi rendszer működését, kivétel a `Bundle-Activator` és az `Import-Package` fejlécek.

A manifest állomány elemei:

Bundle-Name az elkészített bundle emberek számára olvasható neve

Bundle-SymbolicName Java package név, ez azonosítja a bundle-t (az egyetlen kötelező elem)

Bundle-Description a bundle hosszabb szöveges leírása

Bundle-ManifestVersion a bundle által használt OSGi változat verziószáma

Bundle-Version a bundle verziószáma

Bundle-Activator `BundleActivator` interfészt implementáló bundle-ben lévő osztály, mely a bundle telepítése után elindul

Export-Package más bundle-ök számára elérhetővé tett saját csomagok és verziószámaik listája

Import-Package a bundle fordításához és futtatásához szükséges külső csomagok listája

Private-Package olyan saját csomagok, amelyeket nem teszünk elérhetővé más bundle-ök számára

A bundle-ök egy OSGi példányon belül, azonos JVM-ben futnak. Ennek vannak előnyei (teljesítmény növekedés, kisebb erőforráshasználat, interprocessz kommunikációt nem szükséges használni), de hátrányai is (hozzáférési problémák), melyeket az OSGi úgy old meg, hogy minden bundle-höz saját classloader-t rendel.

2.1. Listing. MANIFEST.MF

```
1 Bundle-Name: Hello World
2 Bundle-SymbolicName: org.available.helloworld
3 Bundle-Description: A Hello World bundle
4 Bundle-ManifestVersion: 2
5 Bundle-Version: 1.0.0
6 Bundle-Activator: org.available.helloworld.Activator
7 Export-Package: org.available.helloworld;version="1.0.0"
8 Import-Package: org.osgi.framework;version="1.3.0"
9 Private-Package: org.notavailable.helloworld
```

Ha egy komponens kapcsolatba akar lépni az OSGi keretrendszerben lévő más komponenssel, akkor azt a hozzá tartozó egyedi `BundleContext`-en keresztül teheti meg. A `BundleContext` megszerzéséhez implementálnia kell a `BundleActivator` interfészt (2.1. ábra), melynek `.start()` és `.stop()` metódusainak argumentumaként megkapják a bundle kontextust, illetve előbbi metódusok meghívódnak a komponens indításakor és leállításakor.

2.2. Service, Service Registry

A bundle-ök kiajánlhatnak szolgáltatásokat (*service*), melyekre más bundle-ök feliratkozhatnak. Az OSGi specifikáció szerint a szolgáltatások normál Java objektumok, melyek egy adott interfészt implementálva lettek beregisztrálva az OSGi *Service Registry* moduljába.

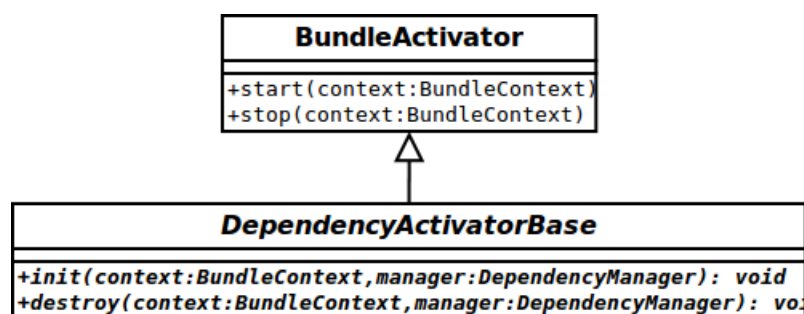
A Service Registry segítségével tudnak tehát a bundle-ök kiajánlani szolgáltatásokat, rajtuk keresztül tudják lekérdezni az elérhető szolgáltatásokat. Ha egy bundle-nek szüksége van egy másik bundle által kiajánlott szolgáltatásra, akkor az OSGi keretrendszeről lekéri a *BundleContext*-en keresztül az adott szolgáltatás referenciáját.

Még robusztusabb használat lehetséges, hogy ha implementáljuk a *ServiceListener* interfészt (*.serviceChanged(ServiceEvent e)* metódust), mellyel értesülhetünk a regisztrált illetve eltávolított szolgáltatásokról, így dinamikusan követhetjük a szolgáltatások elérhetőségét, felkészülhetünk például egy hirtelen kieső szolgáltatás kezelésére. Ugyanezt a funkcionalitást elérhetjük másképpen egy *ServiceTracker* példány használatával is.

A szolgáltatások kezelése ezen módszerek használatával azonban nem skálázódik jól, összetett rendszerekben egyszerűbb technológiák használata jobban kifizetődő. Lássunk tehát néhány példát a szolgáltatások dinamikus kezelésére:

Service Binder A Service Binder komponens létezésének ma már csak történeti okai vannak, mivel a későbbi OSGi R4 specifikációban megjelent újdonságok (Declarative Services, Dependency Manager, iPOJO) teljes mértékben helyettesítik ezt a megoldást. A Service Binder megpróbálja a bonyolult Activator példányok implementációját megkerülni és automatizált módon kezelni a komponens szolgáltatás függőségeit.

A korábbi *BundleActivator* osztályból való leszármaztatás helyett a *GenericActivator* osztályból kell leszármaztatni az osztályunkat, és egy XML fájlban deklaratív módon kell meghatározni, hogy milyen szolgáltatás példányokat akarunk létrehozni és azoknak milyen függőségei vannak.



2.2. ábra. DependencyManager osztálydiagram

Declarative Services A Declarative Services egy komponens model, melynek felépítését befolyásolta a korábbi Service Binder megoldás. A technológia célja, hogy leegyszerűsítse az olyan OSGi komponensek készítését, melyek OSGi szolgáltatásokat használnak.

Fontos jellemzők, hogy használatával nincs szükség explicit módon implementálni a szolgáltatások kiajánlását és felhasználását, ezek deklaratív módon teljes mértékben

XML segítségével végezhető el. Leszármaztatásra vagy interfész implementációra nincs szükség; a szolgáltatások implementációja a *lazy loading* tervezési mintát követi; a komponensek konfigurációja bármikor megtörténhet Configuration Admin szolgáltatás segítségével.

Dependency Manager Ebben az esetben a `DependencyActivatorBase`-ből kell leszármaztatni az osztályunkat, és implementálni az `.init()` és `.destroy()` metódusokat (2.2. ábra). A `DependencyManager` példány segítségével pedig a Decorator mintát használva adhatunk hozzá szolgáltatásokat, és azok függőségeit a komponenshez, valamint konfigurálhatjuk azokat (2.2. kódrészlet).

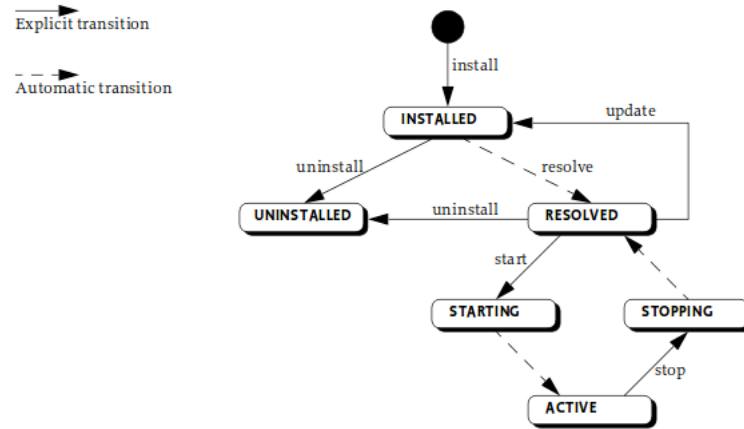
2.2. Listing. Példa a DependencyManager használatára

```
1 manager.add(createService()
2     .setInterface(WebService.class.getName(), null)
3     .setImplementation(WebServiceImpl.class)
4     .add(createServiceDependency()
5         .setService(ConfigurationAdmin.class, null)
6         .setRequired(true))
7     .add(createServiceDependency()
8         .setService(LogService.class, null)
9         .setRequired(false));
```

iPOJO Az iPOJO komponens model központi koncepciója, hogy ne legyen szükség POJO példányoknál bonyolultabb objektumok használatára. Használatát tekintve nagyon hasonló a Declarative Services-hez, azonban az iPOJO kicsit több eszközt kínál a fejlesztőknek. A komponenshez szükséges szolgáltatásokat és azok függőségeit szintén XML fájlban kell deklaratív módon specifikálni, és a rendszer automatikusan kezeli ezeket a függőségeket. Újdonságok a korábbi megoldásokhoz képest például a Java annotációk támogatása, API hozzáférés a belső működéshez, könnyű kiegészíthetőség és beépített távoli konfiguráció JMX segítségével.

2.3. Életciklus

Az OSGi keretrendszer dinamikusságát a *Life-cycle* (életciklus) rendszer szolgáltatja, mely által a bundle-öket futásidőben lehet telepíteni, elindítani, leállítani, frissíteni, eltávolítani más hagyományos alkalmazásokkal ellentétben.



2.3. ábra. Bundle életciklus (forrás: OSGi Service Platform Release 2 [1])

A bundle-ök futtatása előtt megvizsgálja a keretrendszer a csomag futás idejű függőségeit, és ha olyan kielégítetlen függőségeket talál, melyek szükségesek a bundle futtatásához, akkor nem indítja el a komponenst. Egy bundle életciklusának állapotai megfigyelhetők a 2.3. ábrán, valamint az állapotok leírása a 2.1. táblázatban látható.

Állapot neve	Leírás
INSTALLED	A bundle sikeres telepítve lett.
RESOLVED	A bundle-nek minden függősége ki lett elégítve. Kész az elindításra, vagy már le lett állítva.
STARTING	A bundle el lett indítva, de még nincs aktiválva, .start() metódus még nem tért vissza.
ACTIVE	A bundle aktiválva lett és aktív.
STOPPING	A bundle le lett állítva, .stop() metódus még nem tért vissza.
UNINSTALLED	A bundle el lett távolítva, ez az életciklus végállapota.

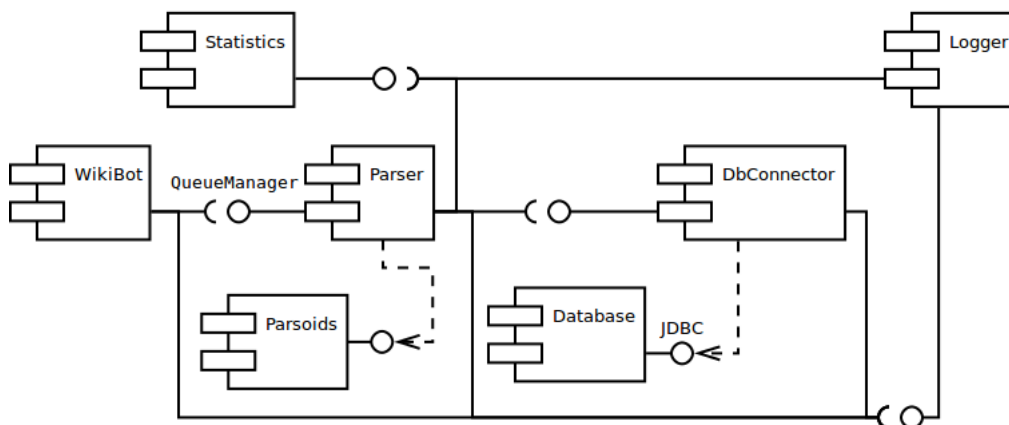
2.1. táblázat. Bundle életciklus állapotai

3. fejezet

Tervezés

Ahogy a 2. fejezetben is láttuk az OSGi keretrendszer megfelel a kért szoftver által támasztott követelményeknek, felhasználásával a rendszer kellően flexibilis lesz. Az 1. fejezetben levont tanulságok alapján érdemes megtervezni a rendszert, illetve az ott említett Wikipédia Miner architektúráját fel lehet használni a tervezés során, az alapvető komponensek meghatározásában segíthet.

Első lépésként meghatároztam a leendő rendszer komponenseit (3.1. ábra), mely képes folyamatosan működve a Wikipédia tartalmát rendszerezett formában eltárolni és azt később a rá épülő alkalmazások számára elérhetővé tenni. Mivel a rendszer első indításakor az adatbázisban még semmilyen adatok nem találhatók, és a Wikipédia frissülő cikkeinek követésével, csak az aktuális változások kerülnek be az adatbázisba, a rendszer indításakor már a Wikipédián korábban közzétett cikkek nem kerülnek be a felépített tudásbázisba. Ezt a problémát úgy lehet a legegyszerűbben megoldani, hogy az on-the-fly feldolgozás mellett a statikus dumpból való adatok importálását is lehetővé tesszük. Ezzel kombináljuk a 1. fejezetben megismert alkalmazások előnyét a követelményekben meghatározottakkal, így egy sokkal erősebb eszköz állítható elő.



3.1. ábra. Az alkalmazás tervezett komponensei

A követelményekből következően szükségszerűen meghatározott komponenseken túl az alkalmazás minőségének javítása érdekében érdemes a komponensek felügyelhetőségét is biztosítani.

Ezáltal az alkalmazás működéséből részletesebb adatok is megfigyelhetővé válnak azon kívül, hogy működik vagy sem a rendszer, illetve a funkcionalitás tesztelése és a teljesítménymérés is sokkal egyszerűbb lesz. A felügyelhetőséget naplózás és különféle metrikák segítségével terveztem lehetővé tenni.

Miután a különálló komponenseket meghatároztam, megkezdődhetett sorban az egyes komponensek megtervezése. Minden komponenst a komponensalapú fejlesztésnek megfelelően egyesével, mint különálló programokat lehetett megtervezni és implementálni. A következőkben bemutatom az egyes komponensek szerepét, alapvető működését és a tervezésük lépéseit.

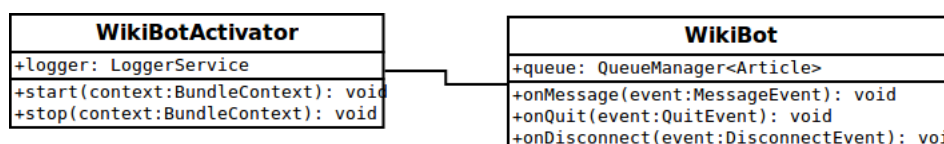
3.1. WikiBot bundle

Az adatok on-the-fly feldolgozásáért felelős ez a komponens. Ezt a tulajdonságot a frissülő cikkek megszerzésének módjával fogom biztosítani. Az adatok könnyű kezelése miatt minden cikknek egy reprezentációját is ki kell alakítani a programban, amelyet tovább kell adni a következő komponensek számára, további feldolgozás céljából. Ezek alapján meghatározott use-case-ek láthatóak a 3.2. ábrán.



3.2. ábra. A WikiBot komponens használati esetei

Megtervezését egy rövidebb előkutatás előzte meg, mely során keresni kellett valamilyen lehetőséget, amelynek segítségével folyamatosan értesülhet az alkalmazás, ha egy új Wikipédia cikk jelenik meg, vagy frissítenek egyet a hivatalos oldalon. A legegyszerűbb megoldásnak végül egy hivatalos, Wikipédia által üzemeltett IRC csatorna tűnt, ahol folyamatosan publikálják az oldalon frissülő cikkeket.



3.3. ábra. A tervezett WikiBot osztálydiagramja

A tervezett komponens tehát egy IRC Bot kliens lesz, amely a fent említett IRC csatornára lesz feliratkozva. A tervezésnél figyelni kellett arra, hogy miután megszerzi az új cikk szükséges adatait, azokat azonnal tovább kell adnia a következő komponensnek. Más teendőket nem végezhet, működését nem szabad hosszú távon feltartani, mert a gyorsan frissülő cikkek miatt elveszhetnek információk, míg a komponens mással foglalkozik.

A komponens két osztályból fog állni, az egyik egy OSGi BundleActivator implementáció (WikiBotActivator osztály), mely menedzseli a komponens indulását, leállítását, a másik pedig a Wikipédia IRC csatornájával való kommunikációt kezeli (WikiBot).

A tervezett működés szerint a rendszer a WikiBotActivator elindulásával kezdődik, melyet az OSGi keretrendszer példányosít számunkra, majd elindítja azt. A lehető leggyorsabb működés úgy érhető el, hogy ne tartsuk fel a komponens működését, ha az új cikk megszerzett adatait eltároljuk egy átmeneti tárolóban (QueueManager osztály). Ez a queue, ahogy a 3.1. ábrán is látható a Parser komponens egy kiejánlott OSGi szolgáltatása. A QueueManager osztály bemutatása a 4.2 alfejezetben folytatódik. A kiejánlott queue referenciáját az OSGi BundleContext-en keresztül lehet majd megszerezni.

Továbbiakban a WikiBot belép a Wikipédia egyik IRC szobájába, ahol az angol nyelvű cikkeket teszik közzé. Az implementált IRC Bot-ban egy eseménykezelőt kell megvalósítani (.onMessage() metódus), ami akkor hívódik meg, ha egy új üzenet érkezik az IRC szobába. Egy példa üzenet figyelhető meg a 3.1. kódrészletben.

3.1. Listing. Példa üzenet az angol nyelvű Wikipédia IRC csatornájából

```
1 (11.16.43) rc-pmtpa: [[History of Vietnam]] http://en.wikipedia
   .org/w/index.php?diff=580135314&oldid=580104406 *
   124.170.231.126 * (+95)
```

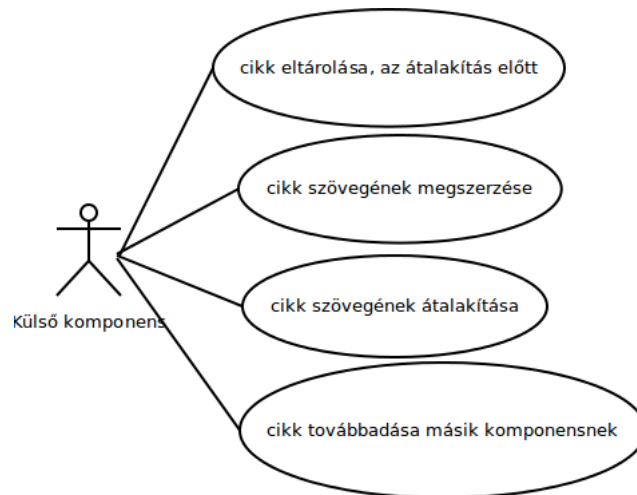
Minden egyes beérkezett üzenet alapján, egy új cikk reprezentációt (Article osztály) hoz létre, mely később végig fog haladni a teljes feldolgozóláncon. A szükséges adatok a cikk címe, mely dupla szögletes zárójelek között található, valamint az új cikk verziószáma, mely az üzenetben található link diff GET paraméterében található. Végül az összeállított cikket a korábban megszerzett QueueManager-ben helyezi el.

Az alapvető követelményeken felül a megfigyelhetőséget is érdemes biztosítani, amelyet a WikiBot komponensnél naplózással oldottam meg. A WikiBotActivator indulásakor a Logger komponens (4.4. alfejezet) referenciáját is megszerzi, melyet naplózásra használhat.

3.2. Parser bundle

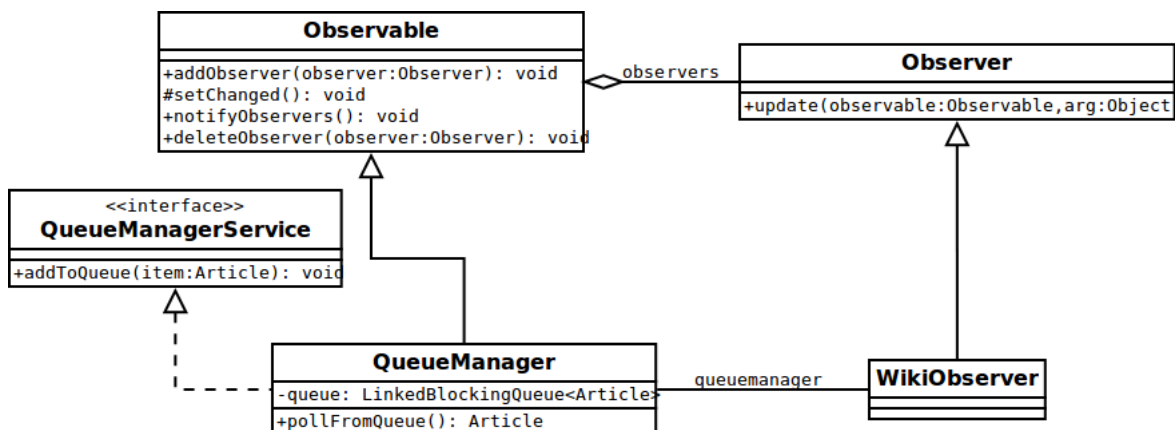
Ebben a komponensben kell lennie valamilyen tároló elemnek, egyfajta queue megoldásnak, melynek szükségessét az előző 4.1. alfejezetben fejtettem ki. Ennek a queue megoldásnak szálbiztosnak kell lennie, hiszen egyszerre fog a WikiBot és a Parser komponens dolgozni vele, ezenkívül a tárolónak blokkolnia is kell beszúrást, ha már túl sok elem van benne. A módosult, vagy újonnan létrehozott cikkek szövegét meg kell szereznie a komponensnek, majd azt át kell alakítania egy meghatározott formátumra. Végül a cikket tovább kell adnia egy másik komponensnek eltárolás céljából. Ezek alapján a use-case diagram elkészíthető (3.4. ábra).

Ahhoz, hogy egy másik komponens el tudjon tárolni a Parser-ben valamit, ahhoz létre kell hozni a Parser komponensben egy szolgáltatást (OSGi service), amelyen keresztül a kommunikáció megtörténhet. Ez az szolgáltatás egyben az egész program mozgatórugója is, hiszen a további use-case-ekben található funkciókat akkor tudjuk végrehajtani, ha legalább egy cikk már el lett tárolva.



3.4. ábra. A Parser komponens használati esetei

Ennél a résznél használtam az *Observer* tervezési mintát [9]: a *QueueManager* osztály *Observable* lett, míg egy *Observer*-ből származó *WikiObserver* kezelte a tárolóba került új elemeket (3.5. ábra).

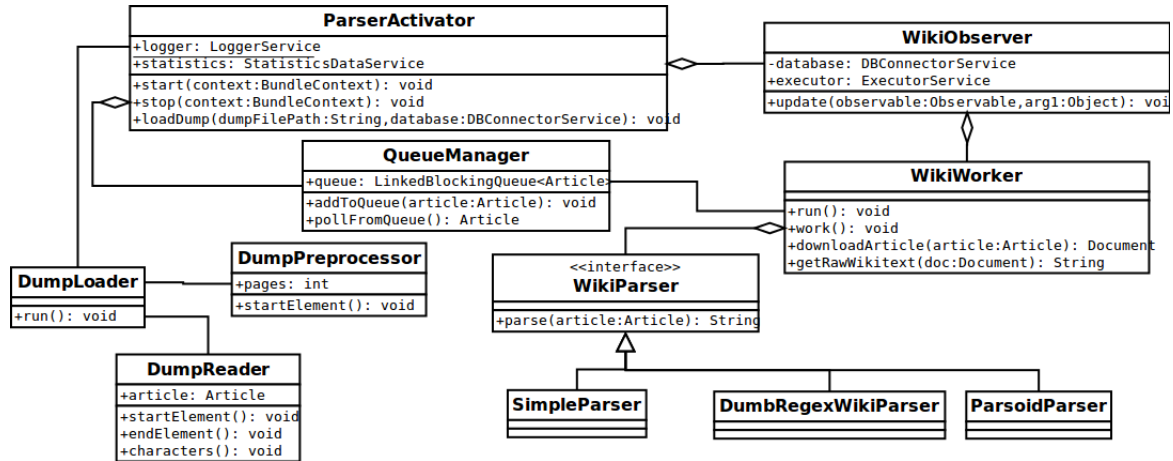


3.5. ábra. Observer minta használata a Parser komponensben

Az Observer tervezési mintának köszönhetően, ha új elem kerül a queue-ba, akkor arról a *WikiObserver* értesülni fog. Mivel a cikkek feldolgozása jelentős időt vehet igénybe, azokat mindenképpen külön szálakon kell megtenni. Ezt a problémát a *ThreadPool* tervezési mintával oldottam meg [6]. A *WikiObserver*-nek (3.6. ábra) van egy *ThreadPool* példánya, és ha értesül a queue frissüléséről, egy új feladatot fog a *ThreadPool*-hoz hozzáadni. A *ThreadPool* mintának megfelelően minden feladat (*WikiWorker*) külön szálon fut.

Ebben a fázisban első lépésben először letöltődnek a hivatalos Wikipédia API-n keresztül a cikkek tartalmai. A Wikipédia a felhasználók számára könnyebben szerkeszthető *Wikitext* formátumban teszi elérhetővé a cikkek tartalmait. A Wikitext egy egyszerű, könnyen használható jelölőnyelv, mely egyértelműen leképezhető a HTML formátumra, így könnyen készíthetők vele webes tartalmak.

Azonban ez a Wikitext formátum kevésbé jól feldolgozható, mint például a HTML formátum, így a cikkek szövegét HTML formába célszerű alakítani. A második fázisban tehát ez az átalakítás



3.6. ábra. A Parser komponens osztálydiagramja

történik meg, ha nem volt még újabb verziójú változat az adott cikkből az adatbázisban. Ehhez ún. parser-eket lehet használni, melyből három felcserélhető, különböző előnyökkel és hátrányokkal rendelkező változat is került a feldolgozóláncba.

- Sztakipedia parser [13]: Könnyen testreszabható és általános célú program, mely a MediaWiki Wikitext formátumról tud HTML formátumra átalakítani. Ez a program egy könnyen kiterjeszthető *Visitor* tervezési mintára alapuló szoftver, mely önálló csomagként is használható a jól felépített API-nak köszönhetően. Használata nagyon egyszerű, használható kimenetet állít elő és sebesség szempontjából is jól teljesít, viszont nagyon nagy méretű fájlknál, például egy teljes Wikipédia dump feldolgozásánál rendkívül sok memóriát használ. Ennek a tulajdonságnak köszönhetően, csak on-the-fly feldolgozás során használható eredményesen ez a parser.
- DumbRegexWikiParser parser: A végletekig leegyszerűsített MediaWiki Wikitext átalakító, mely egy SAX parser implementáció. Az állapotgépes Wikitext feldolgozásnak köszönhetően nem fogyaszt annyi memóriát, mint a Sztakipedia parser, sebesség szempontjából még gyorsabb is az előzőnél, viszont kevésbé használható kimenetet produkál.
- Parsoid parser: Ez az átalakító a MediaWiki által fejlesztett és használt szoftver, mely a Wikitext egy ekvivalens HTML / RDFa kimenetét készíti el, mely automatikus feldolgozásra rendkívül jól használható. Az RDFa szabványban meghatározott attribútumokkal kiegészített HTML kód, sokkal nagyobb jelentéstartalommal bír, mint az alap HTML dokumentum, így a szöveg sokkal jobb lesz további felhasználhatóság szempontjából. A Parsoid egy NodeJS-ben írt szoftver, így felhasználása Java nyelvű alkalmazásokban sokkal nehezebb, mint a korábbi parser megoldásoké, viszont kimenete a célnak leginkább megfelelő, így ez lett az ajánlott átalakító mechanizmus a feldolgozólánban.

Végül a cikk feldolgozásának utolsó, harmadik fázisában eltárolódik a cikk a rendszer adatbázisában. Ebbe a fázisba akkor ér el a feldolgozás, ha nem volt újabb verziójú változat az adott cikkből még az adatbázisban. A verziók meghatározása a cikkek verziószáma alapján

történik. Ha régebbi verziószámú cikk már volt az adatbázisban, akkor azt csak frissíteni kell, ellenkező esetben az adott cikknek még semmilyen változata nem létezik az adatbázisban, így azt újonnan kell beszúrni.

Ennél a komponensnél a megfigyelhetőséget naplózással és különféle metrikák lekérdezhetőségével oldottam meg. A `ParserActivator` indulásakor a `Logger` komponens (4.4. alfejezet), és a `Statistics` komponens (4.5. alfejezet) referenciáját is megszerzi.

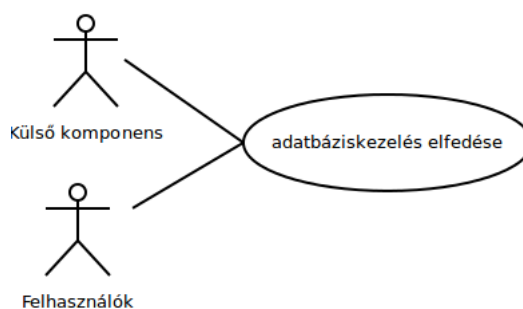
3.2.1. Wikipédia dump feldolgozás

A feldolgozólánc indulásakor, lehetőség van az éppen frissülő cikkek on-the-fly feldolgozásán és adatbázisba való eltárolásán túl a már korábban a Wikipédián elérhető cikkek eltárolására is, melynek forrása a havonta frissülő Wikipédia állománya, amit XML formátumban tesznek elérhetővé.

A dump feldolgozása nagyon hasonlít az on-the-fly feldolgozáshoz, szintén egy `WikiObserver` példány irányítja a feldolgozást és egy `QueueManager`-t használ a cikkek ideiglenes eltárolásához. A hatalmas méretű (2013. szeptemberében 42 Gb) XML fájl, melyben a Wikipédián olvasható cikkek vannak eltárolva, a feldolgozás során először egy előfeldolgozáson megy keresztül, ahol gyorsan néhány statisztikai adatot gyűjt a program az adatbázismentésről. Ilyen statisztikai adat például a dump-ban lévő cikkek száma, így követhető a feldolgozás állapota a tényleges beolvasás során. Ekkora méretű fájlt csak valamilyen állapotgép alapú megoldással lehet feldolgozni, ilyen például *SAX Parser* (Simple API for XML). A parser a megszerzett cikket egy queue-ban tárolja el, ahonnan a `WikiObserver` kiveszi, annak szövegét átalakítja HTML formátumra és eltárolja az adatbázisban.

3.3. Databaseconnector bundle

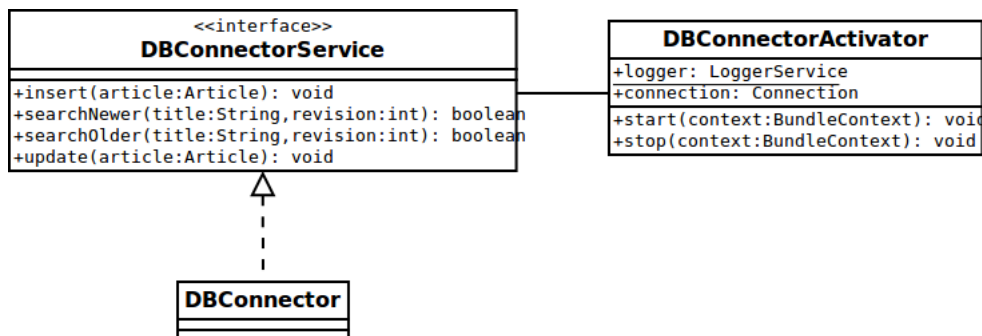
Ebben a komponensben lesz az adatbázisban eltárolva a Parser komponens által feldolgozott minden cikk. Feladatai közé tartoznak tehát az adatbáziskezeléssel kapcsolatos műveletek elfedése más rétegek előtt: cikkek eltárolása, frissítése, törlése és különböző szempontok alapján való keresés a cikkek között (CRUD műveletek). A feldolgozóláncra alapuló további rétegek számára – melyet külső fejlesztők (a rendszer felhasználói) készítenek – ez a komponens egy csatlakozási pont, így kiejánlott szolgáltatását alacsonyabb és magasabb rétegben lévő komponensek is használhatják, az adatbázissal való kommunikáció teljesen el van fedve előlük.



3.7. ábra. A DBConnector használati esetei

A fenti usecase-ek csak általánosan fogalmazzák meg a komponenssel kapcsolatos elvárásokat.

Az egyes feldolgozólánra alapuló komponensek további követelményeket is támaszthatnak a komponens iránt, például különböző szempontok alapján való keresés.



3.8. ábra. A DBConnector komponens osztálydiagramja

A komponens indulásakor (**DBConnectorActivator** osztály `.start()` metódusa), olyan állapotba hozza az adatbázist, hogy az később képessé váljon a fenti feladatokra, valamint egy OSGi szolgáltatást kell kiajánlani, amin keresztül a többi komponens kommunikálhat a **DBConnector**-ral.

3.4. Database

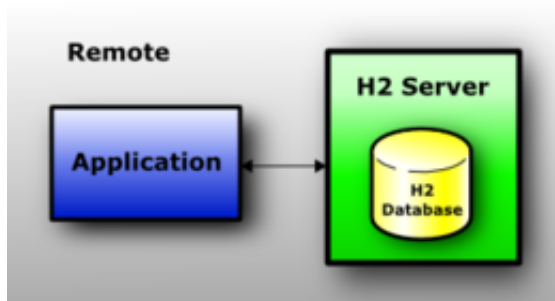
Az adatbázis komponens alapjául szolgáló adatbázis szoftver kiválasztását a WikiBot komponenshez hasonlóan egy rövidebb elő kutatás előzte meg. Fő szempont volt, hogy nagy teljesítményű, gyors és lehetőleg az OSGi keretrendszerrel könnyen integrálható legyen a választott adatbázis megoldás. Két különböző adatbázis lehetőséget vizsgáltam meg, és a tapasztalatok alapján választottam ki a legmegfelelőbbet.

1. Prewayler: Egy nagyon egyszerű és gyors adatbázis, melynek lényege, hogy minden objektumot egyszerű Java objektumokként (POJO) tartunk a memóriában, és a tranzakciókat naplózzuk, az adatbázis pedig a használó programba beágyazottan fut. A tárolandó objektumokra egyetlen megkötés, hogy implementálják a `Serializable` interfészt, és a műveletek is szerializálható objektumok legyenek. A műveleteket a Prewayler naplózza, és időnként teljes mentést (*snapshot*) készít az adatbázisról, így egy rollback esetén a snapshotot visszaállítja és lefuttatja még a szükséges tranzakciókat.

Felépítéséből adódóan sokkal gyorsabb mint bármelyik RDBMS, viszont OSGi környezetben a használata egyelőre nem megoldott. A probléma oka, hogy minden OSGi komponensnek saját classloader-e van, így osztályok dinamikus betöltése, például szerializációval nem lehetséges az OSGi programokban [15]. A megoldás lehetne új `BundleClassLoader` implementálása, azonban ez a Prewayler használatának egyszerűségét rontaná el, így inkább másik megoldást választottam.

2. H2DB: Az összegyűjtött tapasztalatok alapján arra jutottam, hogy OSGi környezetben kevésbé érdemes beágyazott adatbázist választani, ha több komponens között is meg kell oldani a kommunikációt, hiszen belső működésüket tekintve ezen adatbázisok szinte mindig használnak valamilyen szerializációs eljárást.

A H2DB egy objektum relációs adatbázis, mely tud beágyazott és kliens-szerver módban is működni. A készítőik teljesítménystesztei alapján gyorsabb a legtöbb népszerű adatbáziskezelőnél (HSQLDB, Derby, MySQL, PostgreSQL), valamint tervezésekor megpróbálták a HSQLDB és a Derby adatbáziskezelők előnyeit is egyesíteni benne.



3.9. ábra. H2DB szerver módban (forrás: H2 Database Engine [8])

A H2DB további előnyei még, hogy nem szükséges telepíteni, teljes egészében Java nyelven íródott és például Java nyelvű tárolt eljárásokkal is kiegészíthetjük funkcióit. Az eddig felsorolt pozitív tulajdonságok miatt esett a választásom a H2DB-re, bár igaz, hogy beágyazott módban nem használható effektíven, csak úgy, mint a Prevayler. A 3.9. ábrán látható módon, szerver módban is használható a H2DB, így JDBC adatbázishozzáférési API-val lehet kapcsolódni hozzá. A szerver módban való futtatás előnye, hogy sokkal kisebb a csatolás az adatbázissal, így például a DBConnector komponensben, ha a JDBC API-t használjuk, az adatbázis könnyen lecserélhető a rendszerben másik adatbázis megoldásra.

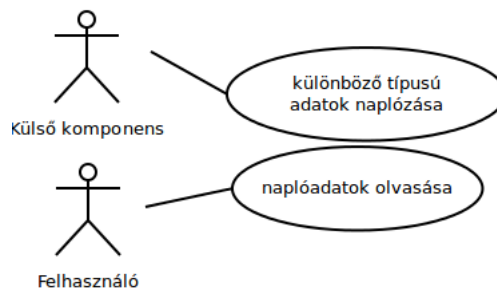
Article	
*title	VARCHAR(255)
°text	CLOB
°length	INT
°revision	INT
°inserted	TIMESTAMP

3.10. ábra. A tervezett adatbázis Article táblája

Az adatbázis kiválasztása után, az adatbázis tábláinak megtervezése következett. Az egyszerű specifikáció miatt egyed kapcsolat diagram rajzolására nincs szükség, hiszen egyetlen egyed van az adatbázisban, mely a cikk reprezentációja (Article). A Wikipediában minden cikket annak címe azonosítja (3.10. ábra), nem lehet két azonos című, de különböző cikk. Így elsődleges kulcs a cím (title) lesz, további szükséges attribútumok a szöveg (text), szöveg hossza (length), szöveg verziószáma (revision), illetve a beillesztés dátuma (inserted).

3.5. Logger bundle

Ha a rendszer állapotait meg akarjuk figyelni, illetve szeretnénk, ha jelezné a hibákat, mindenképpen érdemes az alkalmazást felügyeletre tervezni. Az egyik ilyen lehetőség a naplózás, mely ennek a komponensnek a fő feladata (3.11. ábra).



3.11. ábra. A Logger komponens használati esetei

Első lépésben meg kell tervezni a rendszer felügyeleti modelljét, mely egy állapotgép formájában ábrázolható a rendszer fontosabb állapotaival. Az állapotgép átmenetei események, az átmenetek feltételei hibalehetőségek, a tranzakciók pedig naplózást jelentenek, a használt jelölés:

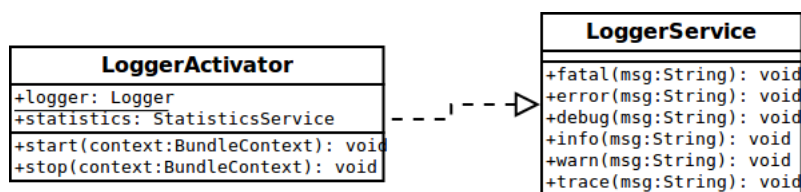
$$\text{esemény}(\text{argumentum})[\text{feltétel}]/\text{tranzakció} \quad (3.1)$$

Naplózás során a rendszer eseményeit több kategóriába soroltam be súlyosság és következmények szempontjából, ezek láthatóak a 3.1. táblázatban.

naplózási szint	rövid leírás
FATAL	Súlyos hiba, amely a rendszer leállításához vezet.
ERROR	Súlyos hiba, azonban a rendszer működése nem feltétlenül áll meg.
WARN	Esetlegesen káros hatással bíró esemény.
INFO	Magasszintű információ a rendszer állapotáról.
DEBUG	Részletes információk a rendszer állapotáról, hibakereséshez használható.
TRACE	Legrészletesebb információk a rendszer állapotairól.

3.1. táblázat. A naplózás szintjei

A fent meghatározott naplózási szinteknek megfelelő naplózási eljárásokat teljesítő komponens egy OSGi szolgáltatás lesz, mely képes a naplózást a követelményeknek megfelelően elvégezni (3.12. ábra).



3.12. ábra. A Logger komponens osztálydiagramja

3.6. Statistics bundle

Az alkalmazás felügyeletének másik lehetőségével meghatározott metrikák lekérdezésére és vizsgálatára van lehetőség. A többi komponens állapotait ennek a komponensnek a segítségével

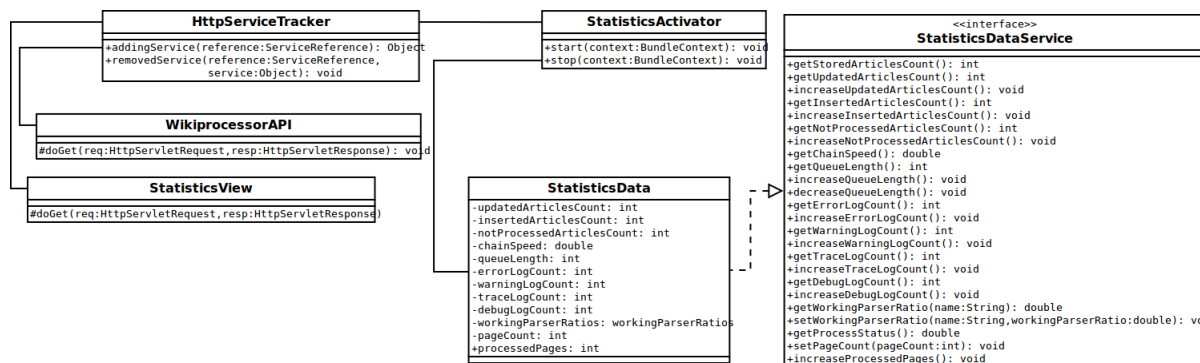
teheti elérhetővé, így a rendszer által végzett feladatról készített statisztikák megfigyelésére is van lehetősége a felhasználóknak (3.13. ábra).



3.13. ábra. A Statistics komponens használati esetei

A feldolgozólánc működését átvizsgálva meghatároztam és kiválasztottam megfigyelés céljából a rendszer állapotának szempontjából fontos mérőszámokat (3.2. táblázat).

Az eddigi követelmények alapján úgy tűnhet, hogy a legjobb választás a JMX (Java Management Extensions) technológia lenne a feladatra, azonban több szempont alapján sem ezt választottam. Bár később a rendszer valamilyen szintű menedzselése szükségessé válhat, mégis OSGi környezetben nem biztos, hogy a JMX a legjobb megoldás. Egyrészt a JMX nem építhető be a már korábban említett OSGi tulajdonságok miatt a rendszerbe, így egy másik implementációt a MOSGi-t (Managed OSGi) kellene használni; másrészt a MOSGi egy sokkal bonyolultabb architektúrát használ és a felhasználók számára sem annyira egyszerű a kezelése.



3.14. ábra. A Statistics komponens osztálydiagramja

A fent felsorolt indokok alapján döntöttem saját menedzselő komponens fejlesztése mellett, mely egy sokkal egyszerűbb architektúrát követ, és nagyon könnyen kiegészíthető. A menedzser komponens egy OSGi szolgáltatásból áll (StatisticsDataService a 3.14. ábrán), melybe a többi komponens rögzítheti az állapotait, illetve két servlet segítségével tud kapcsolatba lépni a felhasználókkal, melyek közül az egyik JSON formátumban elérhetővé teszi a statisztikai adatokat (WikiprocessorAPI), a másik pedig egy HTML oldalt készít a statisztikai adatokkal (StatisticsView). Az elkészített HTML oldalon grafikonok készülnek el dinamikus JavaScript + Ajax segítségével, így az elkészült weboldalon friss statisztikai adatokat lehet megfi-

metrika neve	metrika leírása	metrika mértékegysége	metrika számításának módja
eltárolt cikkek száma	ennyi cikk van az adatbázisban	darab	frissített cikkek száma + újonnan beillesztett cikkek száma
frissített cikkek száma	korábban már az adatbázisban szereplő, frissített cikkek száma	darab	ha kisebb verziószámú azonos című cikk szerepel az adatbázisban, akkor eggyel nő a száma
újonnan beillesztett cikkek száma	korábban még az adatbázisban nem szereplő cikkek száma	darab	ha még nem volt az adatbázisban az adott című cikk akkor eggyel nő a száma
nem feldolgozott cikkek száma	ennyi cikket nem dolgoztunk fel, mert újabb változat van az adatbázisban	darab	minden esetben mikor újabb verziójú cikket dolgoz fel a száma nő eggyel
hibák száma	FATAL, ERROR, WARN, DEBUG, INFO, TRACE üzenetek száma a naplóban	darab	log üzenet írása esetén a statisztikák száma is nő
queue hossza	QueueManager által karbantartott queue-ban lévő elemek száma	darab	QueueManager-be íráskor értéke nő eggyel, elem kivételekor csökken eggyel
feldolgozás sebessége	ennyi cikket dolgoz fel 1 perc alatt a feldolgozólánc	cikk/perc	1 percenként elmentett cikkek számából kivonja az előző elmentett értéket
dolgozóparserek száma	adott nevű observerben ennyi parser dolgozik	százalék	dolgozó parserek száma / az összes elérhető parserek számával
dump feldolgozás állapota	ekkora százaléka lett beolvasva a dump-nak	százalék	beolvasott cikkek száma / összes cikk a dump-ban

3.2. táblázat. A Statistics komponens által megfigyelhető metrikák

gyelni folyamatosan a rendszer állapotáról és a futó folyamatokról.

A Statistics komponens többféleképpen is megszerezheti a feldolgozólánc többi komponensének állapotait. Az egyik lehetőség, hogy saját maga gyűjti az adatokat az eseményekről, a különböző komponensek feliratkoznak a kiajánlott statisztika készítő szolgáltatásra és így frissítik a metrikákra vonatkozó adatokat. Emiatt ezen komponensek (például a Parser) függeni fognak a Statistics komponenstől. Ekkor, ha a parserben lévő queue-ról is szeretnénk adatokat, az egyik lehetőség, hogy a queue minden módosulásánál tájékoztatja a statisztika készítő komponenst a megváltozásról, ez azonban overhead-del jár a queue számára. A másik megoldás, hogy a statisztika készítő komponens kéri le a queue hosszát, ha szüksége van rá, így azonban a Statistics komponens is függeni fog a Parser-től. Emiatt, ha a második lehetőséget választanánk, ami nem jár akkora overhead-del, más a Parser-nél alacsonyabb szintén lévő komponensek nem használhatnák a Statistics komponenst, hiszen a függőségek alapján az OSGi keretrendszer nem tudna előállítani

olyan komponens indítási sorrendet, amelyben minden indítási függőség kielégülne (függőségi hurok alakulna ki). A megmaradt választási lehetőség tehát az, ha minden komponens folyamatosan tájékoztatja a megfigyelhető állapotairól a Statistics komponenst.

4. fejezet

Implementáció

A komponensek fejlesztésekor Eclipse fejlesztő környezetet használtam, verziókezeléshez pedig a Git szoftvert. Az alapértelmezett Eclipse IDE azonban nem tökéletes eszköz az OSGi alapú fejlesztéshez, ezért a Bndtools [4] OSGi fejlesztő keretrendszert Eclipse kiegészítő modulként feltelepítve kezdtem neki a komponensek fejlesztésének.

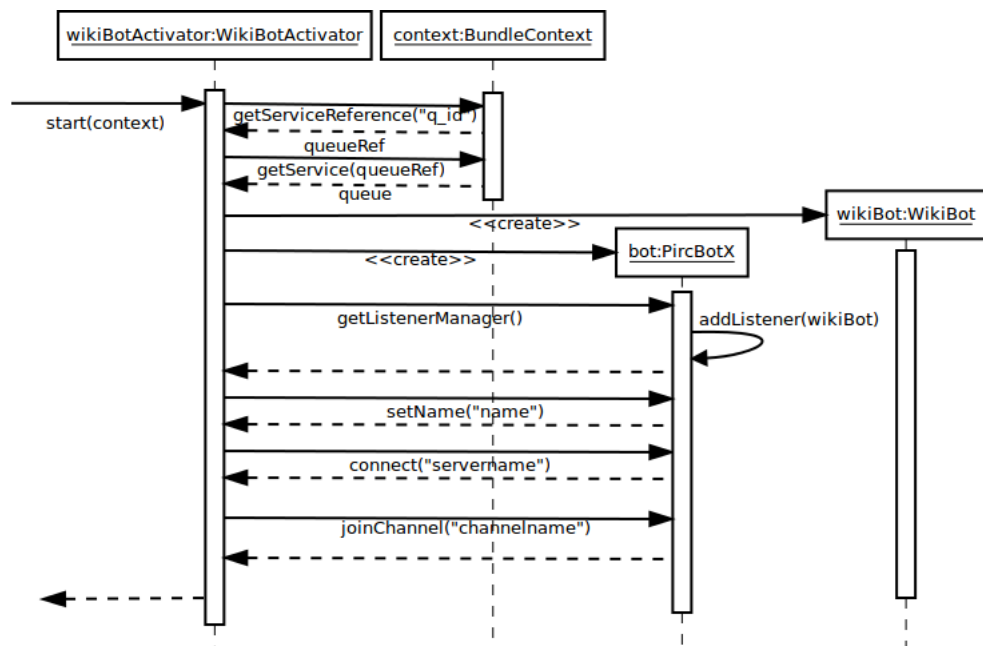
4.1. WikiBot bundle

A WikiBot komponens architektúráját teljes mértékben a fejlesztéskor használt IRC Bot készítő keretrendszerek határozták meg. A kezdeti változatban a PircBot [12] keretrendszert használtam, azonban hosszan tartó használata alatt kiderült, hogy a rendszernek több hibája is van. Hosszabb ideig tartó futáskor ismeretlen eredetű hibák kerültek elő a PircBot könyvtárban, melyek oka-it a rossz tervezés miatt felderíteni sem volt lehetséges. A keretrendszert Java 1.1 verzióra tervezték, nem használja ki az azóta a nyelvben megjelent újdonságokat, így teljesen elavultnak számít technológiailag, hiába fejlesztik azóta folyamatosan. A hibák nagyon sok esetben teljesen elnyelődnek a rendszerben, így sok esetben lehetetlen az alkalmazás hibáinak megfejtése. Egyes metódusok, melyek örökléskor felüldefiniálva jól használhatóak lennének, *private* vagy *final* hozzáférés-vezérlési kulcsszavakkal vannak ellátva. További rossz tervezési döntés a *God Object* antipattern-nel való visszaélés.

Ezen hátrányok miatt egy idő után az újabb, kevésbé elterjedt, de jobban karbantartott verziójára tértem át a PircBot-nak, mely a PircBotX [5] nevet viseli. Ez a váltás azonban csak kis mértékben befolyásolta a fejlesztés menetét, hiszen mindkét keretrendszernél egy esemény alapú architektúra kialakítása szükséges, melyben különféle események esetén implementálni kell a rendszer viselkedését.

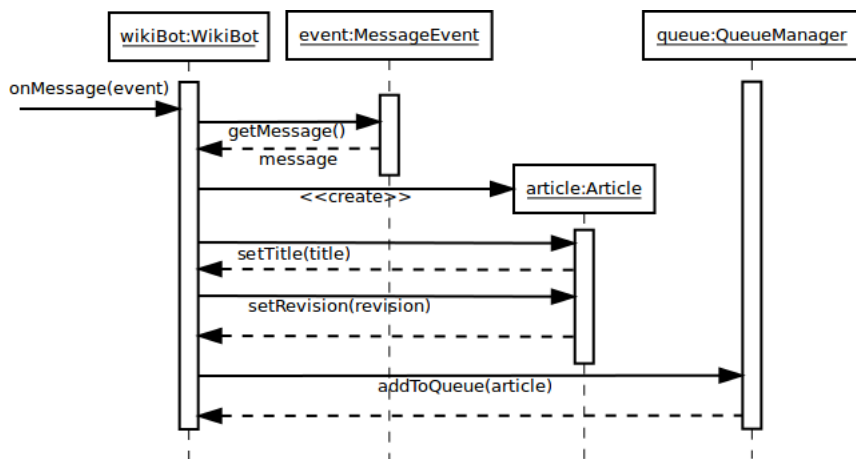
A bundle-ök implementációja során a komponensek függőségeinek, melyek külső könyvtárakban voltak, szintén el kellett készíteni a bundle változatait, hogy azok hiánya ne jelentkezzen hibaként az OSGi függőségkezelőjében. Ez a lépés mind a PircBot, mind a PircBotX esetén a külső könyvtárak OSGi komponens formára való alakítását jelentette, melyet a Bndtools nevű eszközzel hajtottam végre. Ebben a lépésben gyakorlatilag a kész JAR kiterjesztésű library-t kellett egy OSGi specifikáció szerint meghatározott manifest állománnyal újracsomagolni.

A komponens működése látható a 4.1. ábrán: először el kell kérni a `QueueManager` referen-



4.1. ábra. WikiBot komponens indulása

ciáját, majd egy PircBotX példányhoz hozzá kell adni egy `Listener`-t, amely implementálja az eseménykezelőt függvényeket. Ezután meg kell adni az IRCBot nevét, csatlakozni kell egy szervertre, végül be kell lépni egy IRC csatornába. Ettől a ponttól kezdve, ha egy üzenet érkezik a csatornába (4.2. ábra), a WikiBot kiszedi a szükséges információkat reguláris kifejezésekkel és elhelyezi a `QueueManager`-ben.



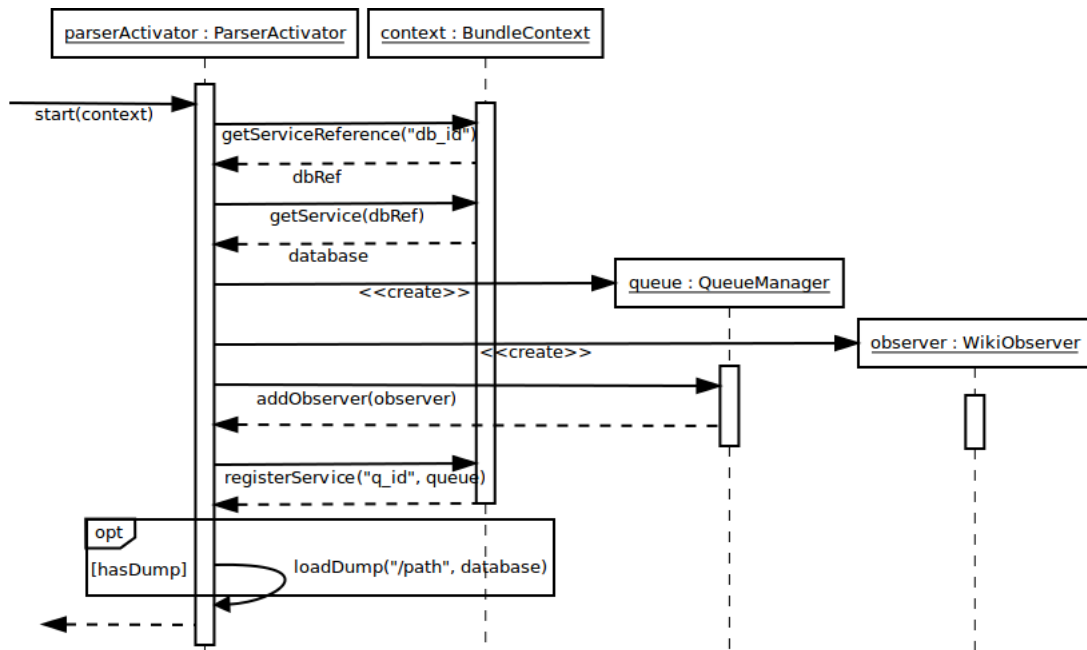
4.2. ábra. Eseménykezelés a WikiBot-ban

4.2. Parser bundle

Az összes komponens közül a Parser felépítése a legösszetettebb, többszálú működése miatt egyszerre történik sok minden, ezért a működésen keresztül fogom bemutatni, hogy teljesíti a Parser bundle a tervezésben meghatározottakat.

A komponens indulásakor először megszerzi a kiejánlott adatbázis szolgáltatás referenciáját, a

feldolgozás után itt fogja eltárolni a feldolgozott cikkeket azok adataival. Létrehoz ezenkívül egy QueueManager-t, beállít (.addObserver() metódus) számára egy újonnan létrehozott megfigyelőt (WikiObserver osztály), és beregisztrálja a queue-t, mint OSGi szolgáltatást. Végül ha elérhető Wikipédia adatbázis mentés (dump), akkor megkezdí annak feldolgozását is. Az indulás folyamata figyelhető meg a 4.3. ábrán.



4.3. ábra. A Parser komponens indulása

A cikkek feldolgozása külön szálakon történik meg, a minél gyorsabb működés érdekében. Minden cikkhez külön Thread objektum példányosítása rendkívül nagy overhead-del járna, ezért a 3 fejezetben már említett ThreadPool mintát használtam. Ez a minta a Thread objektumok újrahasznosításával működik, így nem szükséges mindig új szál példányosítani. Használata a következő 4.1. kódrészletben figyelhető meg, futtatáskor a szálak nevéből látszik, hogy mindig újrahasználjuk a kezdeti 10 darab szálát:

4.1. Listing. Példa a ThreadPool használatára

```

1 public class ThreadPoolExample {
2     public static void main(String args[]) {
3         ExecutorService service = Executors.newFixedThreadPool
4             (10);
5         for (int i = 0; i < 100; i++) {
6             service.execute(new Runnable() {
7                 public void run() {
8                     System.out.println(Thread.currentThread().
9                         getName());
10                    try {
11                        Thread.sleep(1000);
12                    } catch (InterruptedException e) {}
13                }
14            });
15        }
16    }
17 }

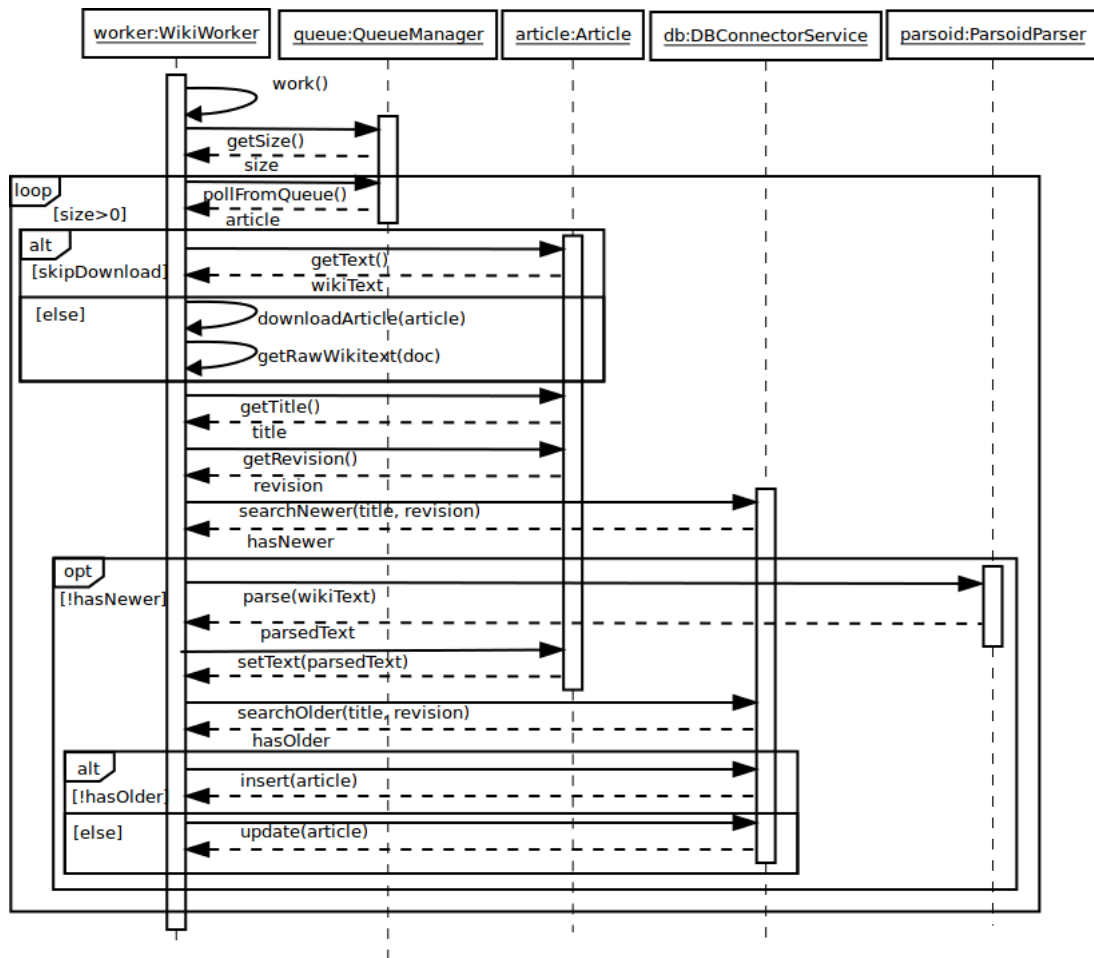
```

```

10         } catch (InterruptedException e) {
11             e.printStackTrace();
12         }
13     }
14     });
15 }
16 }
17 }

```

A Parser működése az eddigieket felhasználva a következő:



4.4. ábra. A WikiWorker szál működése

1. Ha egy új cikk kerül a QueueManager-be, akkor az Observer mintán keresztül a WikiObserver arról értesül.
2. A WikiObserver egy új WikiWorker szálát rendel a ThreadPool-ból a feladathoz, mely megkezdi a cikk előfeldolgozását (4.4. ábrán `.work()` metódus).
3. A szál addig fut, amíg van feldolgozatlan cikk a queue-ban. A dump feldolgozásban fut a szál, akkor nem kell letölteni a cikkek szövegét, mert azokat ilyenkor a dump-ból már

megszereztük. Ellenkező esetben le kell tölteni a hivatalos Wikipédia API-n keresztül a cikket XML formátumban, és ki kell szedni a Wikitext formátumban lévő szövegét.

4. Ha van az adatbázisban már (`.searchNewer()` metódus) újabb verzió az adott cikkből, akkor nem kell tenni semmit, hiszen valószínűleg egy dump feldolgozásban vagyunk és korábban már on-the-fly megszereztük a cikk egy újabb verzióját.
5. Ha nincs újabb verziójú cikk, akkor a Wikitext formátumról HTML (RDFa) formátumra kell alakítani a szöveget (`.parse()` metódus).
6. Ha az adatbázisban van már régebbi verziójú cikk (`.searchOlder()` metódus) az adott cikkből, akkor valószínűleg on-the-fly feldolgozásban vagyunk és korábban már megszereztük a cikk egy korábbi verzióját. Ekkor az adatbázisban csak frissíteni kell a cikk adatait, máskülönben újonnan kell beilleszteni a cikket.

Az előfeldolgozáskor kapott eredmények meghatározzák a feldolgozólánc használhatóságát, így fontos, hogy a megszerzett cikkek Wikitext formátumú szövegét milyen formában tároljuk. A 3. fejezetben láttuk, hogy a feldolgozáskor három cserélhető parser használható a rendszerben, melyek különböző képességekkel bírnak. Ezen parserek a Sztakipedia parser, DumbRegexWikipediaParser és a Parsoid parser. A Sztakipedia parser-nél szükség volt az MTA SZTAKI által készített Java nyelvű könyvtár OSGi komponenssé alakítására.

További felhasználás szempontjából a Parsoid készíti a legjobb kimenetet, azonban egyben a leglassabb is. Ezen lassúság kiküszöbölése miatt különösen jó választás a ThreadPool minta, mert a Parsoid szoftver működése jelentősen függ a bemeneti állományoktól. A Parsoid egy NodeJS-ben írt modul, így a feldolgozólánchoz való illesztését egy szerver oldali NodeJS alkalmazással oldottam meg, mely a HTTP POST kérések törzsében küldött Wikitext szöveget HTML / RDFa formátumra tudja alakítani és azt adja vissza a klienseknek. Ezzel a megoldással a feldolgozólánc HTTP-n keresztül kommunikálhat a Parsoid példányokkal, melyek akár külön szervereken is futhatnak (Parsoids komponens a 3.1. ábrán).

A dump beolvasása szintén ebben a komponensben van implementálva. A Parser indulásakor, ha van elérhető dump file, akkor megkezdí a dump beolvasását. Ezt a folyamatot a `DumpLoader` osztály vezérli, két nagyobb fázisban történik meg a dump beolvasása. A Wikipédia dump-hoz hasonló méretű, óriás XML adatok olvasása semmiképpen nem történhet olyan eszközzel, amiben a teljes adat reprezentációját egyszerre kell tárolni, mint például a DOM Parser. Az én választásom egy állapotgép alapú megoldás lett, a SAX (Simple API for XML) Parser. A feldolgozás első fázisában a tényleges folyamat nyomonkövethetősége érdekében a `DumpPreprocessor` megszámolja a dump-ban lévő cikkek (page csomópontok) számát. A `DumpReader` pedig végigolvassa a dump-ot, és ha beolvasott egy teljes cikket, azt beilleszti egy `QueueManager`-be. Innentől a feldolgozás megegyezik a főszálon történő folyamattal.

4.3. DatabaseConnector bundle

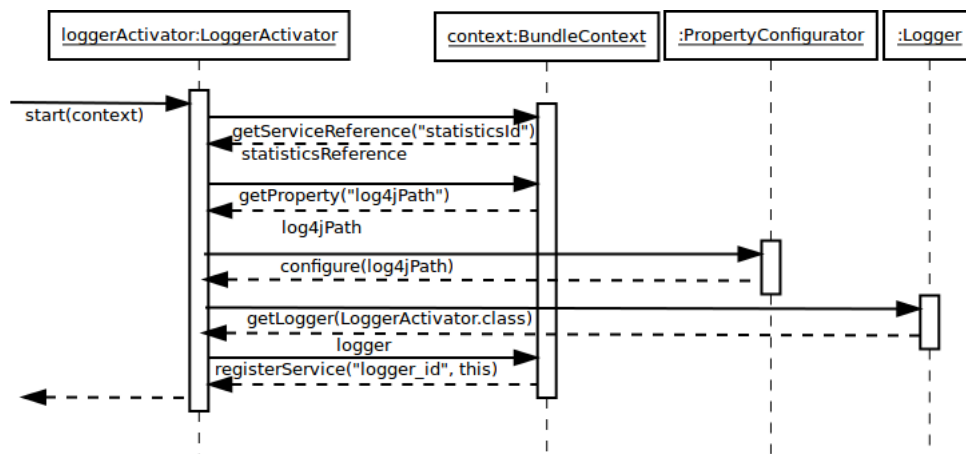
Ez a komponens a tervezésnek megfelelően elfedi az adatbáziskezelést a többi komponensek előtt úgy, hogy felhasználja a kiválasztott H2DB adatbázis JDBC API-ját. A DBConnector kompo-

nens indulásakor először az adatbázis állapotát kell ellenőrizni és szükség szerint olyan állapotba kell hozni, hogy a feldolgozólánc dolgozni tudjon vele. Első lépésben csatlakozni kell a H2DB-hez, majd létre kell hozni hozzá egy adatbázisfelhasználót, egy adatbázist, valamint a szükséges táblákat. Ezzel az adatbázis készen áll a feldolgozólánccal való kommunikációra.

A feldolgozásra kerülő rendkívül nagy adatmennyiség miatt érdemes megfontolni indexek létrehozását is az adatbázisban az Articles táblán. Egy ilyen nagy teljesítményű rendszernél számolni kell az indexhasználat hátrányaival is, ugyanis ha több az írások száma, mint az olvasások száma, az indexek készítése jelentős overhead-et jelenthet az adatbázis számára. Az alkalmazás fő célja, hogy adatokat lehessen kiolvasni belőle, ezért elsődlegesen nem az írások, hanem az olvasások sebességét kell megnövelni. Másrészt az alkalmazás logikája is megköveteli, hogy az olvasás gyors legyen, hiszen ahhoz, hogy egy cikket beillesszünk mindenképpen meg kell nézni, hogy az adott cikk szerepel-e már az adatbázisban. A legtöbbet használt attribútum olvasáskor a cikk címe, így azon hoztam létre indexeket a H2DB-ben.

4.4. Logger bundle

A tervezés részben meghatározott funkcióknak megfelelő módon implementáltam a komponenset. Nem kezdtem saját naplózási rendszer kidolgozásába, hanem a már kiforrott *Apache log4j* megoldást építettem be a komponensbe. Az Apache log4j előnyei, hogy a tervezésénél ügyeltek arra, hogy ne legyen túl nagy hatással az instrumentált kód működésére, és ne kelljen az alkalmazást újrafordítani a naplózás beállításainak módosításához. A log4j beállításait egy konfigurációs fájlban tárolja, ahol széleskörű beállítási lehetőségek használatára van mód.



4.5. ábra. A Logger komponens indulása

OSGi környezetben nem ilyen egyszerű az Apache log4j használata. Alapértelmezett esetben a log4j konfigurációs fájlját (*log4j.properties*) az alkalmazás gyökérkönyvtárába kell helyezni, OSGi környezetben nyilvánvaló módon ez nem működik. A legegyszerűbb megoldás *fragment bundle* készítése lenne, melynek nincs *BundleActivator* osztálya és megosztja classloaderét egy szülő komponenssel. Ennek azonban nagy hátránya van és az Apache log4j egyik előnyét veszítenénk el, mert fragment bundle készítésénél a konfigurációs fájlt a forráskóddal együtt be kellene csomagolni a komponensbe. Így nem lesz konfigurálható a naplózás formája,

szintjei, be lesz égetve egy bundle-be minden beállítás.

A fenti hátrányokat elkerülhetjük, ha *runtime property*-ként megadjuk a helyét a merevlemezén tárolt `log4j.properties` fájlnek a Felix konfigurációjában. A naplózás beállítását egy szöveges fájlban leírhatjuk és azt a rendszerben bárhol elhelyezhetjük. Az Apache log4j beállítását és OSGi szolgáltatásként való kiajánlását a komponens indulásakor kell megtennünk (4.5. ábra).

4.5. Statistics bundle

A Statistics komponens implementálásakor a tervezést követően két servletet és egy adatokat gyűjtő és tároló megoldást kellett létrehozni. A tároló felépítése nagyon egyszerű, a tervezéskor meghatározott metrikákat reprezentáló adatokat tárolja attribútumokban, és implementálja egy kiajánlott OSGi szolgáltatás által meghatározott metrikák módosítására és lekérdezésére szolgáló metódusokat.

A servletek készítése egy kicsit bonyolultabb volt, mert ahhoz, hogy a servlet HTTP-n keresztül elérhető legyen be kell regisztrálni egy servlet container-be. A JEE/J2EE fejlesztés során megszokott `web.xml` fájlban keresztüli konfigurációt nem használhatjuk, OSGi környezetben, ezért az `org.osgi.util.tracker.ServiceTracker` osztályban kell beregisztrálnunk servletünket. Ehhez az előbbi osztályból kellett leszármaztatni egy osztályt (`HttpServiceTracker` a 3.14. ábrán) és felülírni az `.addingService()` és `.removedService()` metódusait, ahol be kellett regisztrálni a servlet osztályokat.

4.2. Listing. A WikiprocessorAPI által szolgáltatott adatok

```
1 { "articlesData": { "eltárolt cikkek száma": 2785, "frissített cikkek  
   száma": 877, "újonnan beillesztett cikkek száma": 1908, "nem  
   feldolgozott cikkek száma": 13 },  
2 "logsData": { "error üzenetek száma": 3, "warning üzenetek száma"  
   : 11, "debug üzenetek száma": 4 },  
3 "queuelength": 6013, "mainworkingparsoids": 0.5714285714285714, "  
   dumpworkingparsoids": 0.8, "dumpprocess": 0.11469880217730487 }
```

A WikiprocessorAPI servlet a metrikák alapján JSON formátumban (4.2) teszi közzé a feldolgozólánc adatait, így ez a servlet gyakorlatilag egy webes API-nak tekinthető. A `StatisticsView` egy statikus HTML oldalt szolgál ki, ahol az előbbi API-t felhasználva JavaScript segítségével diagramokat rajzol ki a servlet. A létrehozott diagramok a metrikákat alakítják a felhasználók számára sokkal jobban megfigyelhető formába.

5. fejezet

Tesztelés és mérések

5.1. A teszteléshez használt eszközök

A feldolgozólánc teszteléséhez elő kellett állítani először a komponensek futtatható változatát. A komponensek implementációjakor használt Bndtools eszköz ebben is nagy segítségemre volt, mert az alkalmazás futtatásakor a Bndtools automatikusan előállította a telepíthető csomagokat. Az egyes komponensek könyvtári függőségeit, amelyek nem részei az alap OSGi futtató környezetnek (tipikusan ezek a harmadik féltől származó könyvtárak) magamnak kellett előállítanom szintén a Bndtools segítségével.

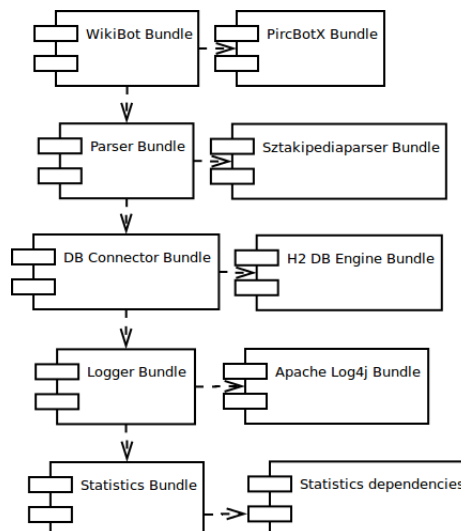
A fordításnál előállított JAR formátumú bundle-ök módosítás nélkül telepíthetők valamilyen OSGi alapú rendszert futtatni képes környezetbe (*bundle repository*-ba). Fejlesztés során a gyorsabb tesztelhetőség miatt az Eclipse saját OSGi implementációs megoldását az Equinox-ot használtam. Egy stabil változat elkészítése után pedig, bizonyos szinten automatizáltan, szkriptek segítségével telepítettem az alkalmazást a SZTAKI Cloud-ban foglalt virtuális gépre, melyet az alkalmazás futtatására készítettem elő tesztelés céljából.

A telepítési környezetben már a könnyen kezelhető Apache Felix [10] OSGi implementációját használtam, mert annak WebConsole eszközével könnyedén kezelhetők a komponensek menedzselése, valamint parancssoros hozzáférés is lehetséges, ahol interaktív módon menedzselhetők a komponensek, sőt még egy egyszerűbb Bash szintaktikáját utánzó szkriptelési lehetőséget is használhatnak a fejlesztők. Teszteléskor a virtuális gépre SSH-n belépve, az Apache Felix konzolos felületét használva telepítettem az elkészített bundle-öket és függőségeiket.

Telepítéshez, illetve a feldolgozólánc komponenseinek frissítéséhez az Apache Felix egy beépített megoldását az Apache Felix Gogo-t használtam, mely egy egységes, szabványos shell felületet határoz meg az OSGi alapú környezetek számára. Amint már említettem az Apache Felix Gogo-val a Unix rendszerekben ismert Bash szintaktikához hasonló parancsokat lehet használni. Ezeket a parancsokat akár szkript formájában is futtathatjuk, ez a Gogo Shell szkript, másnéven *gosh* szkript.

A csomagok függőségeit ábrázolja a 5.1. ábra. A függőségi fa leveleiben található bundle-öket kell először, majd a fában a gyökér felé visszafelé haladva kell a többi bundle-t telepíteni.

A bundle-ök menedzselését a konzolos Apache Felix Gogo Shell-en kívül, egy webes felületen Apache Felix Web Console segítségével is lehet végezni. Segítségével a rendszer állapotát, tulaj-



5.1. ábra. Komponensek függőségei

donságait lehet beállítani és megfigyelni, valamint a bundle-ök részletes tulajdonságait, kapcsolatait lehet felderíteni, illetve beállításait módosítani. A webes felületről a bundle-ök telepítésével, indításával kapcsolatos funkciók szintén elérhetőek.

5.2. A feldolgozólánra épülő kutatómodul

A tesztelés során egy lehetséges felhasználási módot próbáltam ki, mikor a feldolgozólánhoz egy MTA SZTAKI által készített kutatómodult illesztettem. A modul segítségével gépi tanulás, természetes nyelvi feldolgozás témakörökhöz használható program készíthető.

A kutatómodul a feldolgozólán által készített adatbázisból szerzi meg a bemenő adatokat, ehhez a DB Connector komponens kiegészítése volt szükséges, az OSGi szolgáltatást a lekérdezésekhez szükséges metódussal kellett ellátni.

Az kutatómodul az *Annotare* nevet viseli. Célja, hogy XML szabvány szerint jól formázott HTML dokumentumokat nyelvi feldolgozhatóvá és annotálhatóvá tegyen. A könyvtárnak három fő funkciója van:

1. Dokumentum beolvasása SAX Parser segítségével, illetve annak átalakítása. Az átalakítás után keletkezik egyrészt egy plain text dokumentum, amely átadható egy nyelvi feldolgozó modulnak, másrészt egy vagy több lista, amely az XML elemeket, mint annotációkat tartalmazza.

Az Annotare az *UIMA* terminológiáját használja (de az UIMA kódjait vagy könyvtárait egyáltalán nem, azoktól független). Ebben a terminológiában a szóban forgó plain text dokumentum a *Subject of Annotation* (SOFA), az annotációk pedig *FeatureStructure*-ok reprezentálják. A SOFA alapértelmezett konfigurációban az XML elemek tartalmának konkaténálásából áll össze, de az XML tag-ek elhagyásával. A *FeatureStructure*-ok egy begin és end pozícióval rendelkeznek, amely pozíció azt jelöli ki, hogy a SOFA-ban hol kezdődik és végződik egy adott annotáció. A feldolgozáskor alapesetben minden *FeatureStructure* egy XML elemnek felel meg, így a begin és end azt mutatja, hogy hol kezdődik és hol végződik

az adott elem a SOFA stringhez képest. Ez a szám nem azonos azzal a karakterpozícióval, ahol az eredeti dokumentumban az adott XML elem kezdődik, hiszen a SOFA stringbe az XML tag-ek által elfoglalt karakterek nem számítanak bele.

5.1. Listing. Példa az XML SOFA és FeatureStructure kapcsolatára

- 1 XML: <a>xy<c>z</c>
- 2 SOFA: xyz
- 3 FS-ek: (a) 0..11 ; (b) 4..8; (c) 9..11

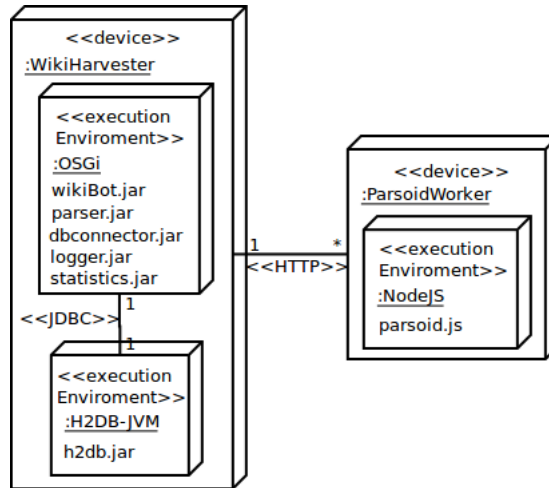
A fentiekén túl egyedi konfigurációval megoldható, hogy az Annotare több, közvetlenül egymás után következő pozíciójú, és azonos típusú FeatureStructure-t egybe olvasson.

2. A második fontos funkció a FeatureStructure lista kiegészítése újabb annotációkkal. Ezek rendszerint egy NLP rendszer (például OpenNLP) kimenetéből származnak. Meg kell adni hozzájuk, hogy milyen XML taggel legyen reprezentálva és hogy mi a begin és end koordinátájuk.
3. Végül az Annotare egy FeatureStructure listát képes XML formátumba szerializálni. Mivel a FeatureStructure-ök átlapolódhatnak, nem kell szigorú tartalmazási viszonyban lenniük (szemben az XML elemekkel), ennek a feladatnak fontos része, hogy szükség esetén a rendszer egy FeatureStructure-t több XML elemre daraboljon fel. Azt, hogy melyik elem lesz több részre bontva az elemek prioritása határozza meg, amely egy konfigurálható paraméter a FeatureStructure-ok esetében.

Az Annotare célja, hogy XML-ként értelmezhető dokumentumokban szereplő szövegeket NLP rendszerekkel feldolgohatóvá tegyen. Tipikusan HTML dokumentumoknál (például a Parsoid parser kimenete) vagy struktúrált dokumentumokat leíró XML fájlok esetén hasznos az alkalmazása. Ezen dokumentumok esetén az a probléma, hogy az XML tag-ek jelenléte teljesen elrontja az NLP rendszerek teljesítményét. Ha ezeket egyszerűen kitöröljük, akkor elveszítünk fontos információkat, ráadásul meg kell oldanunk azt is, hogy az egyik XML elem végén található szó ne olvadjon össze a másik XML elem elején találhatóval. Az Annotare fent részletezett funkciói ebben segítenek. Segítségükkel tetszés szerint vehetünk ki szövegrészeket a dokumentumunkból, feldolgozhatjuk azt, majd az eredményt visszailleszthetjük a struktúrált dokumentumunkba.

5.3. Mérési eredmények

A mérési összeállítás a 5.2. ábrán figyelhető meg. A használt virtuális gépek OpenNebula sablonjának legfontosabb paraméterei a 5.2. kódrészletben olvashatóak. A WikiHarvester gép 4 processzort használ 16 virtuális CPU-mal és 10 Gb RAM-ot érhető el számára. Ezen a gépen fut az OSGi keretrendszer benne a fejlesztett komponensekkel és azok függőségeivel, illetve a H2 adatbázis is itt fut külön folyamatban saját JVM-mel (a H2 adatbázissal JDBC kapcsolaton keresztül történik a kommunikáció). A WikiHarvester gép 30 darab ParsoidWorker virtuális gépet tud használni a cikkek átalakítására a Parser komponensben. A ParsoidWorker gépek 1 darab 2 magos processzort használnak fejenként, 512 Mb RAM memóriával; mindegyikre fel van telepítve a NodeJS, ahol a saját szerveroldali Parsoid megoldás fut.



5.2. ábra. A mérési összeállítás deployment diagramja

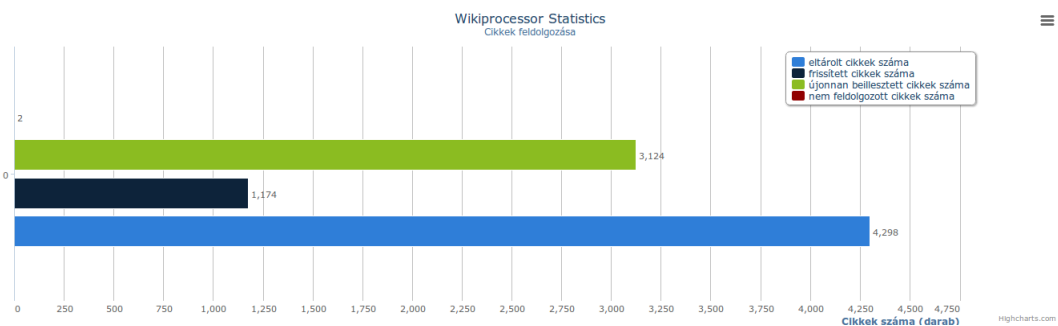
5.2. Listing. Részlet a használt VM-ek sablonjából

```

1 NAME    = [APP]WikiHarvester
2 CPU     = 4
3 VCPU    = 16
4 MEMORY  = 10240
5
6 NAME    = [APP]ParsoidWorker
7 CPU     = 1
8 VCPU    = 2
9 MEMORY  = 512

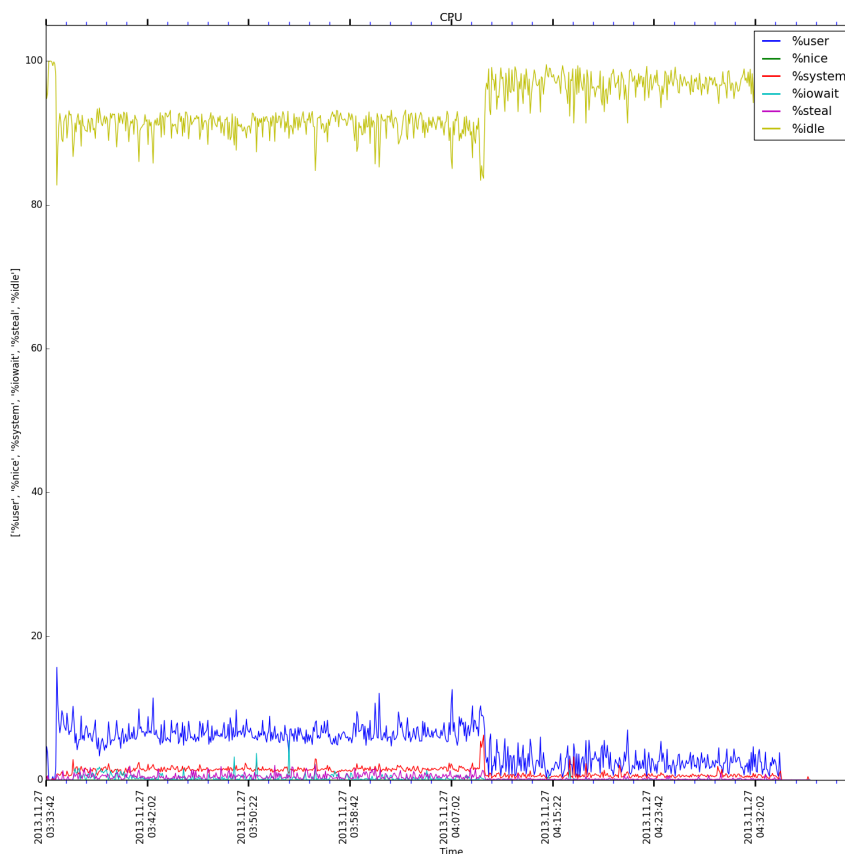
```

A mérés megkezdésekor letöltöttem az angol Wikipédiáról egy adatbázis mentést (dump), melynek mérete 42 Gb volt. A rendszer indításakor megkezdí a feldolgozólánc a dump feldolgozását egy külön szálon, míg a főszálon on-the-fly követi a Wikipédiába kerülő cikkek változásait. A Wikipédia dump feldolgozása egy előfeldolgozó fáissal kezdődik, ahol a későbbi tényleges feldolgozófolyamat követhetőségeért először meg kell számolni, hogy hány cikket kell beilleszteni a dump-ból.



5.3. ábra. A mérés során feldolgozott cikkek statisztikája

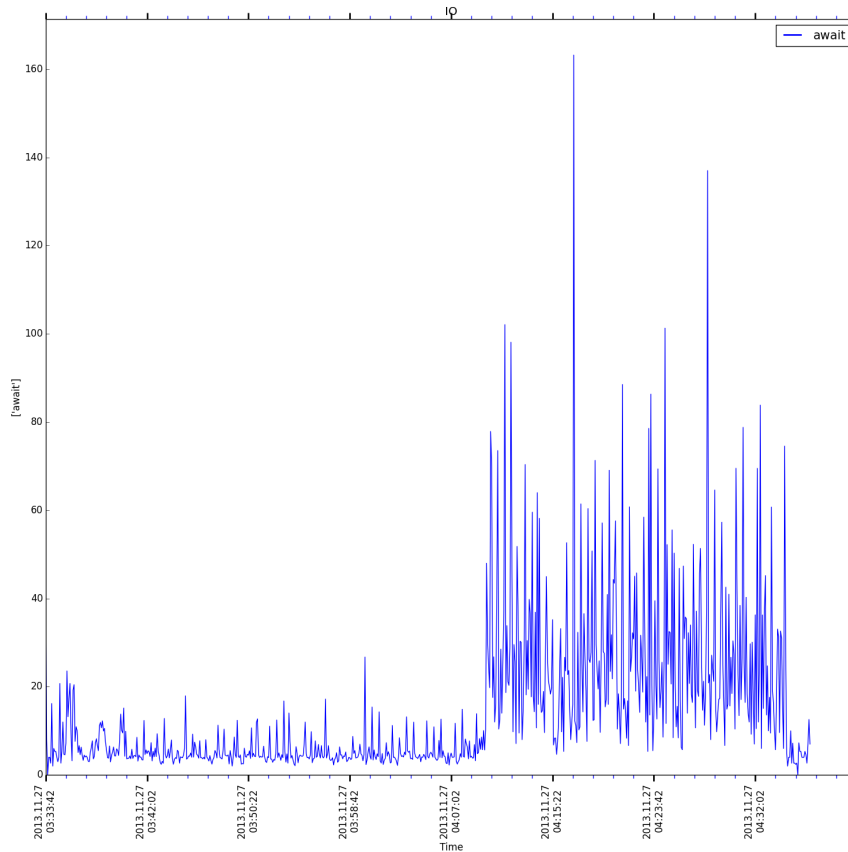
A fent említett előfeldolgozás a mérés során használt Wikipédia mentéssel körülbelül 36 percig tartott, ezalatt 1,36 milliárd cikket olvasott be. Az on-the-fly feldolgozást végző ParsoidWorker-ek száma 25 darab volt, míg a későbbi dump feldolgozást 5 darab ParsoidWorker segítette. Az egy órás mérés végén az adatbázis mérete, üres adatbázissal indítva a mérést 644 Mb lett, ezalatt 4298 cikket dolgozott fel a rendszer. A Statistics komponens által készített részletes feldolgozási statisztikák közül egy látható a 5.3. ábrán.



5.4. ábra. A mérés során mért CPU használat

A teljesítményméréseket az *iostat* nevű ingyenes és nyíltforráskódú monitorozó szoftverrel végeztem, mely képes a processzor és a diszk használatot is monitorozni, valamint a különböző szálak működésének megfigyeléséhez a *top* programot használtam. A processzor használatot ábrázoló diagramon (5.4. ábra) világosan elkülönül az előfeldolgozás fázisa és az a fázis, amikor az on-the-fly és dump feldolgozás egyszerre történik. Az alkalmazás által használt processzoridő (%user, mely a felhasználói módban eltöltött processzoridőt mutatja) körülbelül 2 magot használ ki a rendelkezésre álló erőforrásból, utóbbinál nyilvánvalóan visszaesik kicsit a teljesítmény és átlagosan 1 magnyi erőforrást tud kihasználni a feldolgozólánc.

Kérdés lehet, hogy egy ilyen többszálú alkalmazás, amely több szálát is használ miért nem használja ki a bőséges erőforrást? A válasz az, hogy kihasználja az erőforrásokat, csak a grafikonon ez nem figyelhető közvetlenül meg. Az alkalmazás futásakor a *top* programban világosan látszott, hogy a feldolgozólánc mind a 16 elérhető magot kihasználta a ThreadPool-oknak köszönhetően, azonban a program futásakor legnagyobb részben a processzor más műveletek elvégzésére vár



5.5. ábra. A mérés során mért lemezhasználat

(ezért ilyen magas az `%idle`, azaz tétlenségi idő). Ezen időigényes műveletek hálózati kommunikációt jelentenek: a cikk szövegének letöltése Wikipedia API-n keresztül, kommunikáció a ParsoidWorker-ekkel és várakozás a visszakapott eredményekre. A processzor használat visszaesésére a második részben magyarázat még, hogy a gyorsabb és processzorigényesebb dump feldolgozást végző szálak jobban lefoglalják a processzort a többi szál elől.

A lemezhasználatot bemutató grafikon a rendszer fentiekben meghatározott két állapota még jobban meghatározható, mint a processzorhasználatnál. A lemezhasználati statisztikákból olyan metrikát választottam ki, melyből jó következtetéseket lehet levonni. Ez a metrika az átlagos várakozási idő milliszekundumban, amíg az IO művelet végződik. Látható, hogy az első fázisban viszonylag egyenletesen használja az alkalmazás a diszket, míg a dump feldolgozás fázisában hatalmas kilengések vannak és jobban meghatározza a teljesítményt a hálózati forgalomra való várakozás. Összesítve tehát látható, hogy sem a processzor, sem a diszk teljes erőforrását nem lehet kihasználni a kívánt módon implementált rendszerrel, a hálózati kommunikációra való várakozás miatt; a processzor teljesítménye nem összehasonlítható a hálózati teljesítménnyel.

6. fejezet

Értékelés, továbbfejlesztési lehetőségek

A tesztelés és a mérések során arra jutottam, hogy az elkészített feldolgozólánc képes ellátni azokat a feladatokat, és teljesíti a meghatározott követelményeket, melyek a rendszer tervezése előtt célul lettek kitűzve. Az elkészített alkalmazás úgy lett megtervezve, hogy a korábbi megoldások hiányosságait próbálja pótolni és az on-the-fly feldolgozás miatt hiánypótló is lehet a szemistruktúrált adatforrást feldolgozó rendszerek között.

A hasonló megoldások vizsgálatánál láttuk, hogy egy Wikipédiához hasonló mennyiségű adathalmaz feldolgozása semmilyen esetben sem triviális, a használható feldolgozáshoz idő és számítási kapacitás szükséges. A teszteléskor rendelkezésemre álló erőforrások elegendőek voltak a funkcionalitás ellenőrzéséhez, de számítási kapacitásuk a használhatóságot nem érte el. Korábbi vizsgálatok azonban azt mutatták, hogy a rendszer lineárisan skálázódik, így több erőforrás alkalmazásával a feldolgozólánc minden előnyét ki lehet használni.

A rendszer az OSGi keretrendszer előnyeit teljes mértékben kihasználja, a különböző feladatot ellátó komponensek, a WikiBot, Parser és a DBConnector csatolása rendkívül kicsi és azok könnyen menedzselhetők. Felhasználhatóság szempontjából sokat javított a komponensek funkcionalitásán a felügyeletre tervezés implementálása, a Logger és Statistics komponensek elkészítése.

A SZTAKI által készített elemző kódok integrálása a rendszerbe nagyon könnyű feladat volt. Ez is mutatja, hogy a komponensek jól lettek megtervezve, újabb, más funkcionalitással rendelkező komponens integrálása a rendszerbe sem lenne túl nagy feladat. Remélhetőleg más kutató modulok beépítésére is sor kerül a jövőben, és hasznos segítsége lehet a kutatásoknak a WikiProcessor feldolgozólánc. További kiegészítések és továbbfejlesztési lehetőségek lehetőségeit a következő pontokban foglalom össze.

Parser komponens hálózati művelet nélkül A mérés során láttuk, hogy a rendszer rendkívül sokat vár hálózati kommunikációra. Egy ilyen nagy teljesítményű rendszernél ez hátrányosan befolyásolja a működést, így érdemes lenne a rendszer leglassabb részét a Parsoid parsert beépíteni a feldolgozóláncba, illetve a hálózati kommunikációt minél jobban kikerülni. Ehhez egyrészt a ParsoidWorker-eket a WikiHarvester géppel azonos gépre kell költöztetni, de hálózati kommunikáció megszűntével elért gyorsítás miatt valószínűleg kevesebb példány futtatása is elegendő lenne.

A kommunikáció azonban a NodeJS-ben írt Parsoid és az OSGi platform között nem egy-

szerű feladat, így több alternatívát is érdemes lesz kipróbálni. A különböző technológiát használó alkalmazások kommunikációjára egy lehetséges megoldás lehet az interprocessz kommunikáció (*IPC*), amelyet mindkét nyelvből natív módon is el lehet érni, például *MessageQueue* használatával. Másik lehetőség lehet a NodeJS C++-ban implementált V8 JavaScript motorja és az OSGi közötti kommunikáció *JNI* segítségével. A NodeJS-hez létezik is már JNI wrapper a Java API-hoz, így a kommunikáció megoldható lenne, ennek a modulnak a segítségével.

Feldolgozólánck készítése egyszerűbben, illetve dinamikusan Az eredeti feldolgozólánc architektúrája, melyben nem szerepelt még a dumpfeldolgozás, nagyon egyszerű és átlátható volt, így a célnak megfelelt. A jó tervezésnek köszönhetően a dump feldolgozás is csak két meglévő osztály példányosítását és egy a feladatot végző rész implementációját igényelte. Ebben a helyzetben még mindig elég egyszerű és átlátható a rendszer működése, de látható, hogy az ilyen jellegű kiegészítések szempontjából nem skálázódik egy idő után túl jól a rendszer. Újabb feldolgozási láncok kialakítására a jelenlegiek mellett tehát egyre nehezebbé válik. Ennek a problémának a megoldása lehet, ezen feldolgozólánck létrehozásának általánosítása.

Egyik lehetőség lehet az *Apache Stanbol* nyíltforráskódú és moduláris szoftver stack kiegészítése az eddig elkészített feldolgozólánccal, mely az *Apache Stanbol* részeként használhatná az teljes software stack egyes komponenseit. Az *Apache Stanbol* szintén *Apache Felix*-et használ OSGi futtatókörnyezetként, így ilyen szempontból nem lenne bonyolult a megoldás kivitelezése. A feldolgozólánck dinamikusan előállítását a *Stanbol*-ban elérhető *Enhancement Chain*-ek segítségével válna elérhetővé. Ennek a megoldásnak a hátránya, hogy az *Apache Stanbol* funkcióit nem egészítené ki a feldolgozólánc, hanem inkább csak használna azokból néhányat.

A fent felsorolt okok miatt érdemes lehet, saját implementációt készíteni a feldolgozólánck készítéséhez, könnyen kiegészíthető és továbbfejleszthető architektúrával. Ilyen problémák megoldására készült az *Apache Commons Chain* implementációja, mely a *Chain of Responsibility* és *Command* tervezési minta kombinációja.

A *Chain of Responsibility* minta alapgondolata az, hogy szétválasztja a kérést küldőket a fogadóktól azáltal, hogy több objektumnak adja meg a lehetőséget a kérés lekezelésére. A kérés egy objektumokból álló lánc láncszemeit járja be egészen addig, amíg valamelyik láncszem le nem kezeli.

A *Command* minta parancsok, kérések objektumba ágyazását, és ezen keresztül feladatok delegálását fogalmazza meg. Különbőféle parancsokat rendelhetünk így klienseinkhez (akár dinamikusan is), a parancsokat akár sorba is állíthatjuk, naplózhatjuk, illetve segítségével visszavonható műveleteket kezelhetünk.

Feldolgozólánc perzisztenciája A jelenlegi rendszer leállításával a működéssel kapcsolatos több információ is elveszik azokban a komponensekben, melyek működése nem állapotmentes. Az egyik ilyen kiegészítés lehetne, hogy ha már egyszer feldolgozott a rendszer indulás után egy teljes Wikipédia dump-ot, egy következő indításnál nem érdemes még egyszer megtenni

azt.

Leállításkor lehetőség van az OSGi komponensekben az aktuális futás közben használt információk elmentésére az adatbázisban. Például a QueueManager példányokban tárolt információkat ilyenkor lehetne elmenteni, vagy a Statistics komponens által készített statisztikák eltárolását is ekkor kellene megoldani.

A félév során folyamatosan újabb és újabb technológiákat ismertem meg, ezekben sikerült kellőképpen elmélyedni, így azt hiszem rengeteget tanultam az elvégzett munka alatt. Ezenkívül sikerült egy olyan szoftverrendszer teljes életciklusát végigkövetni és megvalósítani, mely megfelel a kitűzött céloknak és egy használható keretrendszert nyújt azok számára, akik egy ilyen hatalmas tudásbázist akarnak egyszerűen felhasználni, mint a Wikipédia.

Ábrák jegyzéke

1.1. Wikipedia Miner architektúra (forrás: Artificial Intelligence, Wikipedia and Semi-Structured Resources[7])	9
2.1. BundleActivator osztálydiagram	11
2.2. DependencyManager osztálydiagram	13
2.3. Bundle életciklus (forrás: OSGi Service Platform Release 2 [1])	15
3.1. Az alkalmazás tervezett komponensei	16
3.2. A WikiBot komponens használati esetei	17
3.3. A tervezett WikiBot osztálydiagramja	17
3.4. A Parser komponens használati esetei	19
3.5. Observer minta használata a Parser komponensben	19
3.6. A Parser komponens osztálydiagramja	20
3.7. A DBConnector használati esetei	21
3.8. A DBConnector komponens osztálydiagramja	22
3.9. H2DB szerver módban (forrás: H2 Database Engine [8])	23
3.10. A tervezett adatbázis Article táblája	23
3.11. A Logger komponens használati esetei	24
3.12. A Logger komponens osztálydiagramja	24
3.13. A Statistics komponens használati esetei	25
3.14. A Statistics komponens osztálydiagramja	25
4.1. WikiBot komponens indulása	29
4.2. Eseménykezelés a WikiBot-ban	29
4.3. A Parser komponens indulása	30
4.4. A WikiWorker szál működése	31
4.5. A Logger komponens indulása	33
5.1. Komponensek függőségei	36
5.2. A mérési összeállítás deployment diagramja	38
5.3. A mérés során feldolgozott cikkek statisztikája	38
5.4. A mérés során mért CPU használat	39
5.5. A mérés során mért lemezhasználat	40

Táblázatok jegyzéke

2.1. Bundle életciklus állapotai	15
3.1. A naplózás szintjei	24
3.2. A Statistics komponens által megfigyelhető metrikák	26

Forráskódok jegyzéke

2.1.	MANIFEST.MF	12
2.2.	Példa a DependencyManager használatára	14
3.1.	Példa üzenet az angol nyelvű Wikipédia IRC csatornájából	18
4.1.	Példa a ThreadPool használatára	30
4.2.	A WikiprocessorAPI által szolgáltatott adatok	34
5.1.	Példa az XML SOFA és FeatureStructure kapcsolatára	37
5.2.	Részlet a használt VM-ek sablonjából	38

Irodalomjegyzék

- [1] OSGi Alliance. Osgi service platform release 2. *OSGi Alliance Specifications*, pages 1 – 288, 2001.
- [2] Marco Tulio Valente Andre L. C. Tavares. A gentle introduction to osgi. *ACM SIGSOFT Software Engineering Notes*, 33:1 – 5, 2008.
- [3] Peter Murray-Rust CJ Rupp-Advaith Siddharthan Simone Teufel Ben Waldron Ann Copestake, Peter Corbett. An architecture for language processing for scientific texts. *In Proceedings of the 4th UK E-Science All Hands Meeting*, 1:1 – 8, 2006.
- [4] Neil Bartlett. Bndtools Tutorial. <http://bndtools.org/tutorial.html>, 2011. [Online; hozzáférés 2013.04.25].
- [5] Leon Blakey. PircBotX Java IRC Bot. <https://code.google.com/p/pircbotx/>, 2013. [Online; hozzáférés 2013.11.11].
- [6] Joshua Bloch-Joseph Bowbeer David Holmes Doug Lea Brian Göetz, Tim Peierls. *Java Concurrency In Practice*. Addison Wesley Professional, 2006.
- [7] Simone Paolo Ponzetto Eduard Hovy, Roberto Navigli. Artificial intelligence, wikipedia and semi-structured resources. *Artificial Intelligence*, 194:1 – 252, 2013.
- [8] H2 Database Engine. H2 Database Engine. <http://www.h2database.com/>, 2013. [Online; hozzáférés 2013.11.12].
- [9] Ralph Johnson John Vlissides Erich Gamma, Richard Helm. *Design Patterns: Elements of Reusable Object-Oriented Software*. Kiskapu Kft., 2004.
- [10] Richard S. Hall. Apache Felix OSGi Tutorial. <http://felix.apache.org/site/apache-felix-osgi-tutorial.html>, 2011. [Online; hozzáférés 2013.09.25].
- [11] Klaus Berberich Edwin Lewis-Kelham Gerard de Melo Gerhard Weikum Johannes Hoffart, Fabian M. Suchanek. Yago2: Exploring and querying world knowledge in time, space, context, and many languages. *Proceeding WWW '11 Proceedings of the 20th international conference companion on World wide web*, pages 229 – 232, 2011.
- [12] Paul Mutton. PircBot Java IRC Bot. <http://www.jibble.org/pircbot.php>, 2013. [Online; hozzáférés 2013.11.11].

- [13] Solt Illés Farkas Tamás Oláh Tibor, Héder Mihály. Sztakipédia Parser. <https://code.google.com/p/sztakipedia-parser/>, 2012. [Online; hozzáférés 2013.11.11].
- [14] Humberto Cervantes Richard S. Hall. Challenges in building service-oriented applications for osgi. *Communications Magazine, IEEE*, 42:144 – 149, 2004.
- [15] Humberto Cervantes Richard S. Hall. An OSGi implementation and experience report. In *Proceedings of the Consumer Communications and Networking Conference (CCNC 2004)*, pages 394 – 399, 2004.
- [16] Simone Paolo Ponzetto Roberto Navigli. BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artificial Intelligence*, 193:217 – 250, 2012.
- [17] Da Qing Zhang Tao Gu, Hung Keng Pung. Toward an osgi-based infrastructure for context-aware applications. *IEEE Pervasive Computing*, 3:66–74, 2004.
- [18] Iryna Gurevych Torsten Zesch, Christof Müller. Extracting lexical semantic knowledge from wikipedia and wiktory. *Proceedings of the 6th International Conference on Language Resources and Evaluation*, pages 1 – 7, 2008.
- [19] Benjamin Borschinger Cacilia Zirn Anas Elghafari Vivi Nastase, Michael Strube. Wikinet: A very large scale multi-lingual concept network. *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, pages 1 – 8, 2010.

Függelék

6.1. WikiProcessor

A feldolgozólánchoz tartozó kódok Apache License 2.0 licenz alá tartoznak. A forráskódok, és a működtetéshez szükséges segédprogramok és szkriptek szintén ebben a repository-ban találhatóak meg.



Webes elérés:

<https://bitbucket.org/thesnapdragon/wikiprocessor>

Git URL:

git@bitbucket.org:thesnapdragon/wikiprocessor.git

Licensz:

Apache License, Version 2.0