

# **UNIVERSIDAD EAFIT**

## **Proyecto Final (Segunda Parte)**



**DANIEL MÚNERA SÁNCHEZ**

**JAVIER QUINTERO GONZÁLEZ**

**TÓPICOS ESPECIALES EN TELEMÁTICA**

**SEGUNDA ENTREGA PROYECTO FINAL**  
**TÓPICOS ESPECIALES EN TELEMÁTICA**

**PRESENTADO POR:**  
**DANIEL MÚNERA SÁNCHEZ**  
**JAVIER QUINTERO GONZÁLEZ**

**PRESENTADO A:**  
**EDWIN MONTOYA MÚNERA**  
**PROFESOR**

**UNIVERSIDAD EAFIT**  
**MEDELLÍN**  
**2012**

**1. DESCRIPCIÓN DE LA APLICACIÓN O DEL SERVICIO A  
IMPLEMENTAR.**

## MODELADO Y DISEÑO GLOBAL DE LA APLICACIÓN

NUESTRA IDEA ES REALIZAR UN TIPO DE ARAÑA DE BÚSQUEDA MÁS ESPECÍFICA, CONOCIDA TAMBIÉN COMO TOPICAL CRAWLER, A LA CUAL LE DAREMOS UNA SERIE DE SITIOS WEB DE MULTIPLEX QUE OPERAN EN EL PAÍS, EN LOS CUALES CONSULTAREMOS INFORMACIÓN RELEVANTE SOBRE LAS PELÍCULAS DISPONIBLES, SUS CARTELERAS Y LUGARES PARA VERLAS DEPENDIENDO DE LA CIUDAD.

EL OBJETIVO PRINCIPAL QUE TENEMOS ES FACILITARLE A LOS USUARIOS LA BÚSQUEDA DE PELÍCULAS DISPONIBLES Y HORARIOS, CREANDO UN LUGAR DONDE SE CENTRALICE LA INFORMACIÓN, EVITANDO TENER QUE ENTRAR A 3 O MÁS SITIOS PARA VER QUE PELICULAS ESTAN DISPONIBLES.

PENSANDO EN IMPLEMENTACIÓN, CONSIDERAMOS DOS ESTRATEGIAS:

- O IMPLEMENTAR UN CRAWLER UTILIZANDO LIBRERÍAS COMO ANEMONE:**  
ESTA ESTRATEGIA FUE MUY LLAMATIVA Y NOS LLEVÓ A INVESTIGAR SOBRE CRAWLER Y TOPICAL CRAWLERS, PERO LLEGAMOS A LA CONCLUSIÓN DE QUE TENIENDO LAS PÁGINAS SOBRE MULTIPLEX DE LA CIUDAD IDENTIFICADAS (QUE UTILIZAMOS COMO SEEDS DE BÚSQUEDA), ERA MÁS EFICIENTE REALIZAR BÚSQUEDAS ESTÁTICAS Y PREOCUPARNOS POR ENCONTRAR LA MANERA DE REALIZAR EL PARSING DE LA URL, PARA OBTENER LA INFORMACIÓN REQUERIDA, YA QUE ENCONTRAMOS ALGUNAS DIFICULTADES PARA REALIZARLO.
  
- O IMPLEMENTAR UNA BÚSQUEDA MÁS ESPECÍFICA, UTILIZANDO PRINCIPIOS DE LA BÚSQUEDA POR MEDIO DE TOPICAL CRAWLERS:** LLEGAMOS A ESTA DECISIÓN PENSANDO EN EL PERFORMANCE Y HACIENDO UN BALANCE ENTRE BENEFICIO Y COSTO DE IMPLEMENTACIÓN. TAMBIÉN NOS BASAMOS EN UN EJEMPLO REAL DE UNA EMPRESA BRASILEÑA QUE REALIZA ALGO SIMILAR A NUESTRO TRABAJO Y UTILIZANDO ESTA ESTRATEGIA. ([HTTP://SIEVE.COM.BR/](http://sieve.com.br/))

EXPLICACIÓN DE LA ESTRATEGIA A UTILIZAR:

HACIENDO UNA ANÁLISIS DE LA INFORMACIÓN QUE TENEMOS QUE RECUPERAR, ENCONTRAMOS LOS SIGUIENTES DATOS:

1. CIUDADES DONDE ESTÁN LOS MULTIPLEX
2. LUGARES
3. PELÍCULAS
4. HORARIOS

CARACTERÍSTICAS DE LOS SITIOS WEB DE LOS MULTIPLEX QUE VAMOS A UTILIZAR PARA OBTENER LA INFORMACIÓN:

PARA ESTA PRÁCTICA QUEREMOS EMPEZAR EN 2 SITIOS:

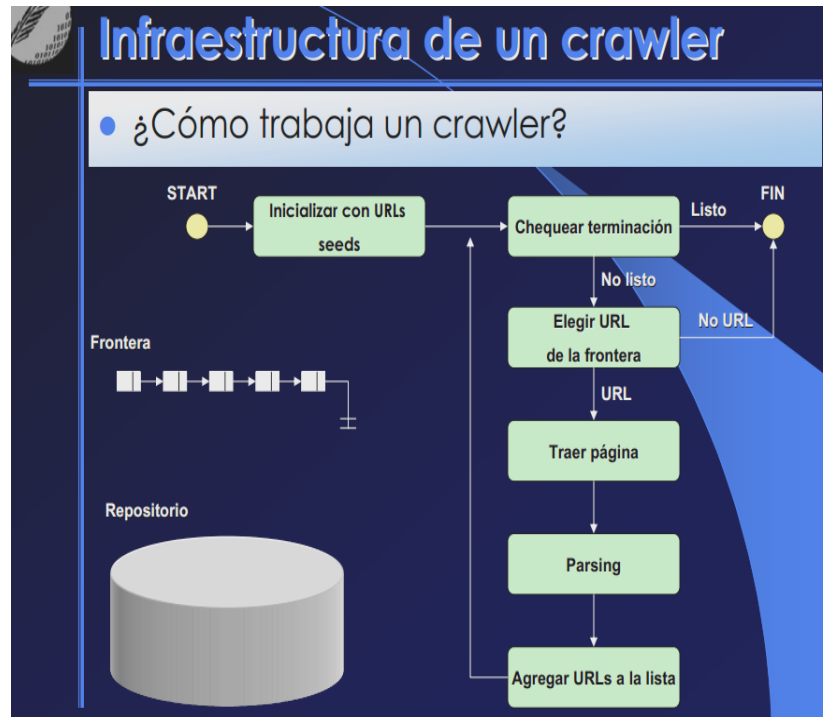
- CINE COLOMBIA
- PROCINAL

ESTOS SITIOS WEB NO TIENEN NINGÚN WEB SERVICES EL CUAL PUDIÉRAMOS APROVECHAR PARA OBTENER LA INFORMACIÓN REQUERIDA, POR TAL MOTIVO PENSAMOS EN UTILIZAR ESTOS

SITIOS WEB COMO SEMILLAS DE UNA BÚSQUEDA Y DESDE AHÍ IDENTIFICAR EN QUÉ PARTE DEL SITIO WEB SE ENCONTRABA LA INFORMACIÓN REQUERIDA. EL PROCESO CONSISTE:

1. TENER IDENTIFICADAS LAS URL'S RELEVANTES PARA NOSOTROS.
2. HACER UN REQUEST PIDIENDO EL HTML DE LA PÁGINA WEB.
3. HACER UN PARSING Y OBTENER LA INFORMACIÓN NECESARIA.

CON EL SIGUIENTE MODELO EXPLICAREMOS QUÉ CONCEPTOS DE TOPICAL CRAWLER UTILIZAMOS PARA CREAR NUESTRA ESTRATEGIA:



AQUÍ SE MUESTRA EL PROCESO GENERAL DE UN CRAWLER, INICIA CON UNOS SEEDS, DESDE LOS CUALES EMPIEZA A NAVEGAR Y A BUSCAR ENLACES DENTRO DEL SITIO Y MIENTRAS SEA NECESARIO SIGUE BUSCANDO Y VA GUARDANDO LOS SITIOS QUE SON RELEVANTES PARA LA BÚSQUEDA. LUEGO DE ESE PROCEDIMIENTO SE PROCEDE AL PARSING PARA GUARDAR LA INFORMACIÓN RELEVANTE.

EN ESENCIA ESE ES EL MISMO PROCEDIMIENTO QUE NOSOTROS VAMOS A SEGUIR PARA OBTENER LA INFORMACIÓN, SOLO QUE CON UN CONJUNTO DE SEEDS MUY REDUCIDOS Y CON PARSERS ESPECÍFICOS PARA CADA SEED UTILIZADA.

#### MANEJO DE LA INFORMACIÓN Y SEGMENTACIÓN DE LA INFORMACIÓN A BUSCAR:

LUEGO REALIZAR UN ANÁLISIS DEL PROBLEMA IDENTIFICAMOS DOS TIPOS DE INFORMACIÓN A BUSCAR:

1. **INFORMACIÓN ESTÁTICA:** ES LA INFORMACIÓN QUE ES NECESARIA PARA REALIZAR LA BÚSQUEDA CON ÉXITO QUE NO ESTÁ SUJETA A CAMBIOS EN EL CORTO PLAZO. ESTOS DATOS SON:

- O CIUDADES
- O LUGARES

ES UN TRABAJO DE MUCHO TIEMPO PODER MONTAR UN MULTIPLEX EN UNA CIUDAD NUEVA, POR ESE MOTIVO CREEMOS QUE NO ES NECESARIO RECOLECTAR ESTA

INFORMACIÓN CADA VEZ QUE SE REALICE UNA BÚSQUEDA. PENSAMOS DARLE LA OPCIÓN A UN ADMINISTRADOR DEL SITIO DE ACTUALIZAR ESTA INFORMACIÓN, O CREAR UN PROCESO QUE SE EJECUTE CADA MES O UNA VENTANA DE TIEMPO SIMILAR PARA ACTUALIZAR ESA INFORMACIÓN

**2. INFORMACIÓN DINÁMICA:** ES LA INFORMACIÓN QUE ESTÁ CAMBIANDO CONSTANTEMENTE Y NECESITA SER OBTENIDA CADA VEZ QUE SE REALICE UNA BÚSQUEDA. ESTOS DATOS SON:

- O PELÍCULAS DISPONIBLES
- O HORARIOS

TENIENDO YA ESTA CLASIFICACIÓN DE LA INFORMACIÓN DECIDIMOS, CREAR UNA BASE DE DATOS PARA GUARDAR LA INFORMACIÓN ESTÁTICA Y AYUDAR A VOLVER LA BÚSQUEDA MÁS RÁPIDA, YA QUE SON DOS CICLOS DE TRABAJO IDÉNTICOS AL DEL TOPICAL CRAWLER QUE ELIMINAMOS. ADEMÁS LA CARGA DE LA INFORMACIÓN GENERAL PARA PODER INICIAR LA BÚSQUEDA TAMBIÉN SE CARGARÁ MÁS RÁPIDO, YA QUE LAS CIUDADES Y LOS MULTIPLEX SE CARGARAN LOCALMENTE EN LOS SELECTORES AL INICIAR LA APLICACIÓN.

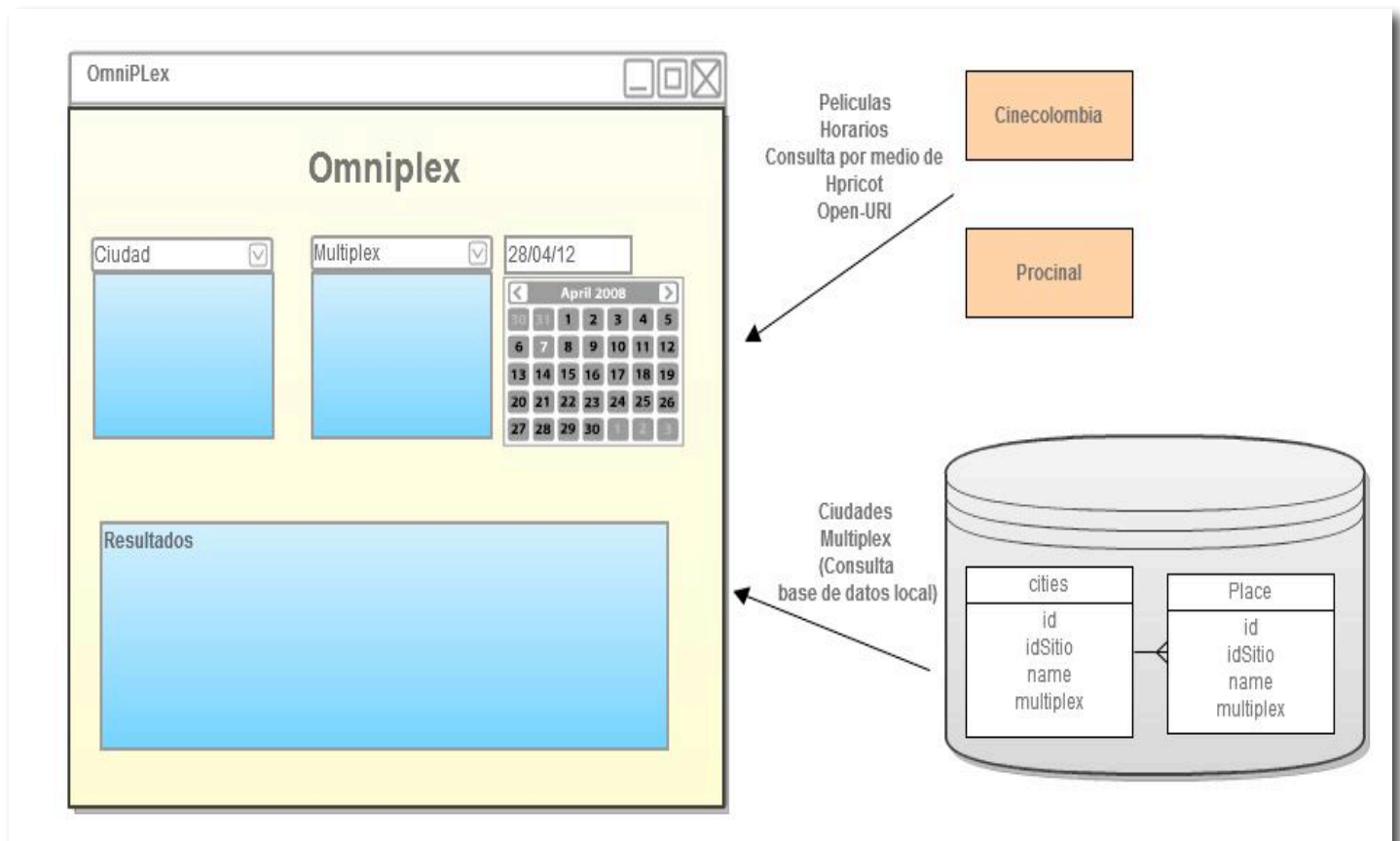
#### **DETALLES A TENER EN CUENTA:**

ALGUNOS DETALLES QUE SON CRÍTICOS PARA NOSOTROS SON:

- **EL MANEJO E INTEGRIDAD DE LA INFORMACIÓN:** COMO SABEMOS, CADA PÁGINA TIENE SU FORMA DE PRESENTAR LA INFORMACIÓN, POR EJEMPLO, EL NOMBRE DE LAS CIUDADES PUEDE ESTAR EN ENCODING UTF-8 U OTRO, O PUEDE TENER CIERTAS DIFERENCIAS ENTRE MAYÚSCULAS Y MINÚSCULAS. POR LO QUE TENEMOS QUE BUSCAR ALGORITMOS PARA ENCONTRAR LA DISTANCIA ENTRE STRINGS E IDENTIFICAR CUANDO DOS CIUDADES SON DIFERENTES Y CUANDO NO. Y ASÍ UNIFICAR EL LISTADO DE CIUDADES DE TODOS LOS MULTIPLEX QUE SE VAN AGREGANDO AL SISTEMA.
- EL DISEÑO DE LOS PARSERS PUEDE SER COMPLICADO, DEBIDO A LA ESTRUCTURA DE ALGUNOS SITIOS QUE EN MUCHAS OCASIONES TIENEN CONTENIDO CARGADO DINÁMICAMENTE CON JAVASCRIPT O CONTENIDO EN FLASH. EN ALGUNOS CASOS LAS PÁGINAS SON UNA COMBINACIÓN DE SITIOS WEB Y POR TAL LAS CARTELERAS DE LAS DIFERENTES CIUDADES TIENEN CAMINOS MUY DIFERENTES PARA SER ENCONTRADAS.
- EL TIEMPO DE BÚSQUEDA TIENE QUE SER CORTO Y EFECTIVO PARA QUE LA APLICACIÓN PUEDA SER EXITOSA.

#### **MODELO Y DISEÑO GLOBAL**

TENIENDO UN POCO MÁS CLAROS LOS DETALLES Y EL ESTADO DEL ARTE DE NUESTRO PROBLEMA, LLEGAMOS AL SIGUIENTE MODELO PARA MOSTRAR LA LÓGICA DE INTERACCIÓN.



## 2. MARCO DE REFERENCIA Y ESTADO DEL ARTE DEL MODELO EMERGENTE

LA ENORME CANTIDAD DE INFORMACIÓN DISPONIBLE EN LA WEB HA HECHO NECESARIO EL DESARROLLO DE RECURSOS QUE PERMITAN ENCONTRAR DE MANERA EFICIENTE Y PRECISA LA INFORMACIÓN REQUERIDA POR EL USUARIO. DENTRO DE LAS DIVERSAS TÉCNICAS DESARROLLADAS CRAWLING ES CAPAZ DE RASTREAR DETERMINADOS TÓPICOS DE MANERA RÁPIDA SIN TENER QUE EXPLORAR TODAS LAS PÁGINAS WEB.

MIRANDO LA IMPORTANCIA DE ESTE TIPO DE HERRAMIENTAS DECIDIMOS DESARROLLAR UNA **APLICACIÓN WEB BASADA EN LOS CONCEPTOS DE CRAWLING**, ESPECIFICAMENTE DEL CONCEPTO DE TOPICAL CRAWLING. EL SISTEMA SOLO CONSISTE EN UNA APLICACIÓN WEB, QUE TIENE LA CAPACIDAD RECOPIRAR INFORMACIÓN DISTRIBUIDA EN LA WEB Y EXPONERLA AL USUARIO PARA SU FACIL ACCESO. A CONTINUACIÓN EXPONEMOS UNA EXPLICACIÓN DE LOS CONCEPTOS EN LOS CUALES NOS BASAMOS PARA SUSTENTAR NUESTRO TRABAJO.

TOPIC -SPECIFIC WEB CRAWLER HAN SIDO DESARROLLADOS CON EL OBJETIVO DE RECOPIAR PÁGINAS DE INTERNET ELIGIENDO SÓLO LAS RELACIONADAS CON UN TEMA EN PARTICULAR. ALGUNOS CRAWLER DESARROLLADOS SON:

- **WTMS:** WEB TOPIC MANAGEMENT SYSTEM. PARA RECOLECTAR INFORMACIÓN RELACIONADA CON UN TEMA EN PARTICULAR EL USUARIO TIENE QUE ESPECIFICAR ALGUNAS URLS RELEVANTES QUE SE CONVERTIRÁN EN LA BASE PARA LAS BÚSQUEDAS. EL WTMS DESCARGA LAS URLS BASE Y CREA UN DOCUMENTO VECTOR REPRESENTATIVO (RDV) BASADO EN LA FRECUENCIA DE OCURRENCIA DE LAS PALABRAS CLAVES EN ESAS PÁGINAS. A CONTINUACIÓN SE DESCARGAN LAS PÁGINAS QUE SON REFERENCIADAS POR LAS URLS SEMILLA Y SE CALCULA SU SIMILITUD CON EL RDV QUE HA SIDO DEFINIDO PREVIAMENTE Y SE PASA UN RANGO ESTABLECIDO SE AÑADEN A UNA COLA. SE CONTINÚA DE ESTA MANERA HASTA QUE LA COLA ESTÁ LLENA O SE ESTABLECE UN LÍMITE. DESPUÉS DE RASTREO, LA COLECCIÓN SE COMPONE DE LAS URLS SEMILLA, ASÍ COMO TODAS LAS PÁGINAS SIMILARES A ESTAS URLS DE SEMILLAS QUE TIENEN RUTAS HACIA O DESDE LAS SEMILLAS.

ESTA COLECCIÓN ES UNA BUENA FUENTE DE INFORMACIÓN DISPONIBLE EN LA WEB PARA EL TEMA ESPECIFICADO POR EL USUARIO. EL RASTREADOR TIENE UNA LISTA DE DIRECCIONES URL DE PARADA PARA EVITAR LA DESCARGA DE ALGUNAS PÁGINAS POPULARES (COMO YAHOO Y LAS PÁGINAS DE NETSCAPE), COMO PARTE DE LA COLECCIÓN.

- **SHARK-SEARCH:** SE BASA EN EL ALGORITMO FISH-SEARCH EN EL CUAL LA CLAVE PRINCIPAL ES LA SIGUIENTE: SE TOMA COMO ENTRADA UNA URL SEMILLA Y UNA COLA DE BÚSQUEDA Y DINÁMICAMENTE SE CREA UNA LISTA DE PRIORIDADES (INICIALIZADA CON LAS URL SEMILLA) DE LAS SIGUIENTES URLS (LLAMADAS NODOS) SER EXPLORADAS. EN CADA PASO EL PRIMER NODO SE EXTRAE DE LA LISTA Y ES PROCESADO. CADA TEXTO ES ANALIZADO Y EVALUADO DE SI ES RELEVANTE O NO Y SE DECIDE SI SE CONTINÚA EXPLORANDO EN ESTA DIRECCIÓN O NO. CUANDO UN DOCUMENTO PROCEDENTE DE UNA URL ES RECUPERADO SE ANALIZA EN BUSCA DE LINKS. LOS NODOS PROCEDENTES DE ESTE DOCUMENTO SON LLAMADOS HIJOS Y LE ES ASIGNADO UN VALOR DE PROFUNDIDAD. SI EL PADRE ES RELEVANTE, LA PROFUNDIDAD DEL HIJO ES UN CONJUNTO DE VALORES PREDEFINIDOS, DE OTRA MANERA, LA PROFUNDIDAD DE LOS HIJOS ES UN CONJUNTO QUE ES UN GRADO MENOR QUE EL DEL PADRE. CUANDO LA PROFUNDIDAD LLEGA A CERO, LA DIRECCIÓN SE ELIMINA Y NINGUNO DE LOS HIJOS ES INSERTADO EN LA LISTA.

ALGUNAS MEJORAS SON HECHAS A ESTE ALGORITMO Y AHORA SE UTILIZA UNA PUNTUACIÓN "DIFUSA" DE RELEVANCIA, DANDO AL HIJO UNA PUNTUACIÓN HEREDADA, POR LO QUE PREFIEREN LOS HIJOS DE UN NODO QUE TIENE UN MEJOR RESULTADO. UNA MEJORA MÁS SIGNIFICATIVA CONSISTE EN REFINAR EL CÁLCULO DE LA PUNTUACIÓN POTENCIAL DE LOS HIJOS NO SÓLO MEDIANTE LA PROPAGACIÓN DE LAS PUNTUACIONES MÁS RELEVANTES DE LOS ANCESTROS DE RELEVANCIA EN LA JERARQUÍA, SINO TAMBIÉN HACIENDO USO DE LA META-INFORMACIÓN CONTENIDA EN LOS ENLACES A LOS DOCUMENTOS.

- Get as Input parameters, the initial node, the width (*width*), depth (*D*) and size (*S*) of the desired graph to be explored, the time limit, and a search query
  - Set the depth of the initial node as *depth* = *D*, and Insert it into an empty list
  - While the list is not empty, and the number of processed nodes is less than *S*, and the time limit is not reached
    1. Pop the first node from the list and make it the current node
    2. Compute the relevance of the current node
    3. If *depth* > 0:
      1. If *current\_node* is irrelevant
        - Then
          - For each child, *child\_node*, of the first *width* children of *current\_node*
            - Set *potential\_score*(*child\_node*) = 0.5
          - For each child, *child\_node*, of the rest of the children of *current\_node*
            - Set *potential\_score*(*child\_node*) = 0
        - Else
          - For each child, *child\_node*, of the first ( $\alpha * \text{width}$ ) children of *current\_node* (where  $\alpha$  is a pre-defined constant typically set to 1.5)
            - Set *potential\_score*(*child\_node*) = 1
          - For each child, *child\_node*, of the rest of the children of *current\_node*
            - Set *potential\_score*(*child\_node*) = 0
      2. For each child, *child\_node*, of *current\_node*,
        - If *child\_node* already exists in the priority list,
          - Then
            1. Compute the maximum between the existing score in the list to the newly computed potential score
            2. Replace the existing score in the list by that maximum
            3. Move *child\_node* to its correct location in the sorted list if necessary
        - Else Insert *child\_node* at its right location in the sorted list according to its *potential\_score* value
3. For each child, *child\_node*, of *current\_node*,
  - Compute its depth, *depth*(*child\_node*), as follows:
    1. If *current\_node* is relevant,
      - Then Set *depth*(*child\_node*) = *D*
      - Else *depth*(*child\_node*) = *depth*(*current\_node*) - 1
    2. If *child\_node* already exists in the priority list
      - Then
        1. Compute the maximum between the existing depth in the list to the newly computed depth
        2. Replace the existing depth in the list by that maximum.
- EndWhile

---

**Figure 1: The Fish-Search Algorithm**

---

- **INFOSPIDER:** BUSCA ONLINE INFORMACIÓN RELEVANTE PARA EL USUARIO, TOMANDO DECISIONES AUTÓNOMAS SOBRE QUÉ ENLACES SEGUIR. UNA PARTE IMPORTANTE DEL SISTEMA ES QUE EL USUARIO PUEDE REGENERAR LA IMPORTANCIA DE UN DOCUMENTO VISITADO POR INFOSPIDERS, ESTO OCURRE DE FORMA ASÍNCRONA CON RESPECTO A LA BÚSQUEDA EN LÍNEA, Y ALTERAR EL COMPORTAMIENTO SUBSECUENTE DE LOS AGENTES, CAMBIANDO EL AMBIENTE DE BÚSQUEDA. ES DECIR, EL USUARIO OBRA RECÍPROCAMENTE CON EL AMBIENTE PARA PREDISPONER EL PROCESO DE LA BÚSQUEDA. EL USUARIO PROPORCIONA INICIALMENTE UNA LISTA DE PALABRAS CLAVES Y UNA LISTA DE LOS PUNTOS DE PARTIDA, EN FORMA DE UN ARCHIVO DE BOOKMARK. CADA AGENTE ES ASIGNADO A UN ENLACE SOBRE EL QUE ACTÚA SEGÚN SU REPRESENTACIÓN.

UNA FASE DE GRAN IMPORTANCIA DENTRO DEL PROCESO REALIZADO POR EL CRAWLER ES LA FASE DE CONOCIMIENTO EN LA CUAL SE DESTACA LA CAPACIDAD DE APRENDIZAJE DEL RASTREADOR EL CUAL NECESITARÁ COMO ENTRADA POR PARTE DEL USUARIO ALGUNAS



PALABRAS CLAVES, O DEPENDIENDO DEL TIPO DE CRAWLER DESARROLLADO, LAS URLS BASE. ESTE PROCESO CONSTA DE TRES FASES:

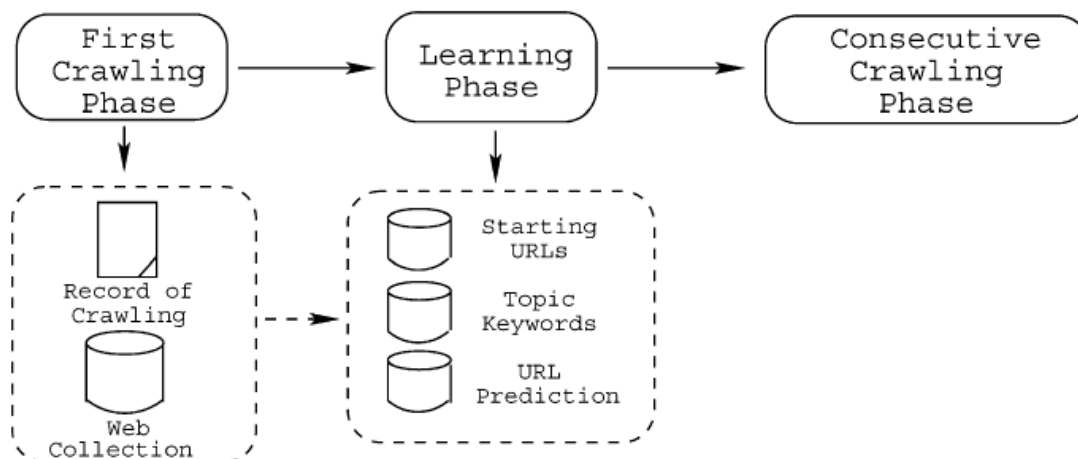


Fig. 1. The diagram of learnable crawling process comprises of three phases: first crawling, learning and consecutive crawling.

- **FIRST CRAWLING PHASE:** DURANTE ESTA FASE EL RASTREADOR SÓLO TIENE PALABRAS CLAVE QUE DESCRIBEN UN TEMA DE INTERÉS SUMINISTRADO POR EL USUARIO. EL RASTREADOR A CONTINUACIÓN ENVIA LAS PALABRAS CLAVES A UN MOTOR DE BÚSQUEDA PARA CONSTRUIR EL CONJUNTO CANDIDATO DE DIRECCIONES DE PARTIDA, ESTA ES UNA OPCIÓN, EN LA OTRA EL USUARIO DA COMO PUNTO DE PARTIDA ALGUNAS DIRECCIONES URL. LOS RESULTADOS DE ESTA FASE SON UNA COLECCIÓN DE PÁGINAS WEB Y UN RASTREO DE REGISTROS. EL RASTREADOR UTILIZA ESTOS RESULTADOS PARA APRENDER Y CONSTRUIR BASES DE CONOCIMIENTO EN LA FASE DE APRENDIZAJE.
- **LEARNING PHASE:** ES UNA FASE EN EL QUE EL RASTREADOR INTENTA APRENDER CÓMO RECOLECTAR DE LA MEJOR MANERA LA INFORMACIÓN TENIENDO COMO BASE LAS EXPERIENCIAS PREVIAS. LA INFORMACIÓN RECOLECTADA SERÁ CATEGORIZADA Y CLASIFICADA SEGÚN LOS ESTÁNDARES PROPIOS Y SE ALMACENARÁN LOS QUE TENGAN LA MEJOR PUNTUACIÓN.
- **CONSECUTIVE CRAWLING PHASE:** FINALMENTE EL RASTREADOR RECOPILA PÁGINAS WEB MEDIANTE EL USO DE LAS BASES DE CONOCIMIENTO CREADAS QUE SERVIRÁ COMO BASE PARA NUEVAS BÚSQUEDAS.

## ESTRUCTURA DE WEB CRAWLER

LA ESTRUCTURA DE UN WEB CRAWLER ES LA SIGUIENTE:

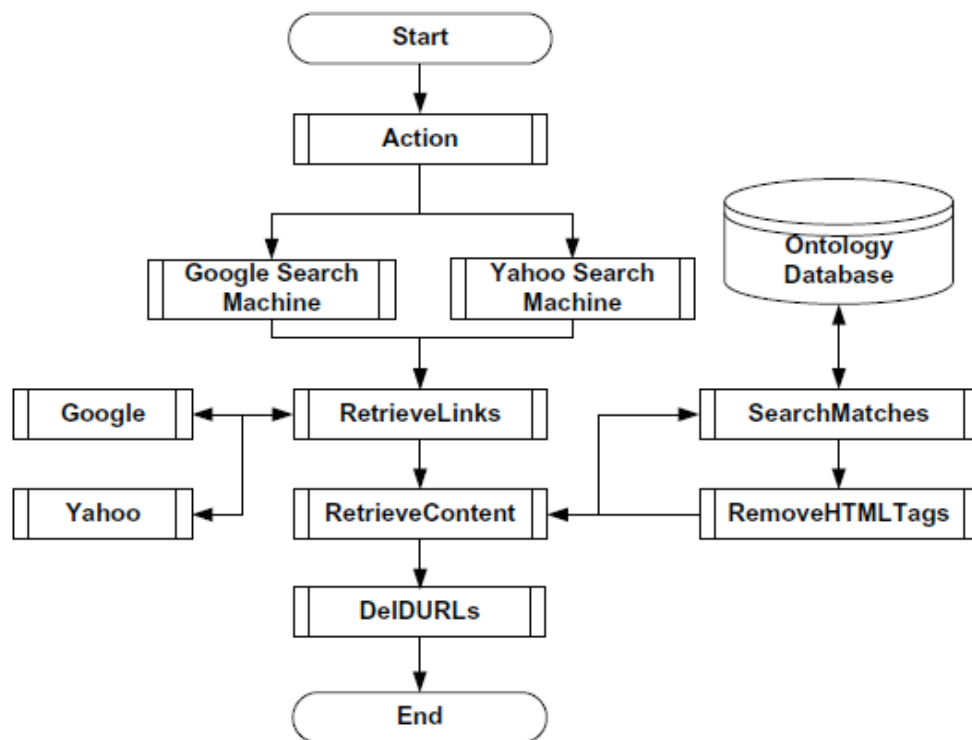


Fig. 4. System structure of Web Crawler.

- **ACTION:** TRANSFORMA UNA CONSULTA INTERNA EN UNA URI Y ES EMBEDIDO DENTRO DE UNA CONSULTA DE GOOGLE O YAHOO SEGÚN SEA EL CASO.
- **GOOGLE/YAHOO SEARCH MACHINE:** DECLARA UN OBJETO URL Y AGREGA LA URL DE GOOGLE EN LA CONSULTA DEL CÓDIGO URI Y ENTONCES SE UTILIZA UN BUCLE INTERATIVO PARA LEER SU CONTENIDO LÍNEA POR LÍNEA. POR ÚLTIMO, MUESTRA EL CONTENIDO COMO ARCHIVO DE TEXTO PARA EL ANÁLISIS FINAL.
- **RETRIEVE LINKS:** USA EXPRESIONES REGULARES PARA BUSCAR LA INFORMACIÓN QUE HA SIDO DESCARGADA DE LA URL. RETORNA LOS HIPERVÍNCULOS Y LA SALIDA DE ELLOS EN UN ARCHIVO DE TEXTO PARA PROVEER AL SISTEMA PARA NUEVOS PROCESAMIENTOS.
- **RETRIEVE CONTENT:** UTILIZA EL ARCHIVO DE ENLACE PARA LEER SU CONTENIDO LÍNEA POR LÍNEA
- **SEARCH MATCHES:** APOYO INTERNO A LA FUNCIÓN ANTERIOR PARA SABER SI LA PÁGINA ANTERIOR ESTÁ DENTRO DEL RANGO QUE NOS INTERESA CONSULTAR. DE SER ASÍ SE AÑADE A LA BASE DE DATOS EN LA CUAL ESTAMOS GUARDANDO LA INFORMACIÓN RELEVANTE PARA NUESTRA CONSULTA.

- **REMOVE HTML TAGS:** FUNCIÓN DE APOYO A LOS DOS ANTERIORES Y SE ENCARGA DE ELIMINAR LAS ETIQUETAS HTML DE LAS PÁGINAS QUE HAN SIDO SELECCIONADAS.
- **DELDURLS:** CONTINÚA CON LA COMPARACIÓN DE LAS DIRECCIONES OBTENIDAS POR LOS MOTORES DE BÚSQUEDA PARA EVITAR LA DUPLICACIÓN DE PÁGINAS.

EN TÉRMINOS GENERALES ESTÁ ES LA FORMA EN LA QUE OPERA UN WEB CRAWLER, SE OBTIENEN PÁGINAS DE INICIO SOBRE LAS CUÁLES SE OPERARÁ, DICHAS PÁGINAS PUEDEN SER ASIGNADAS DE MANERA ESTÁTICA (COMO ES NUESTRO CASO) O SE CONSIGUEN A TRAVÉS DE MOTORES DE BÚSQUEDA COMO GOOGLE O YAHOO. POSTERIORMENTE SE PROCEDE A ANALIZAR LA INFORMACIÓN QUE SE HA ADQUIRIDO, ES AQUÍ DONDE ES IMPORTANTE EL TRABAJO DEL ALGORITMO DE BÚSQUEDA QUE SE ENCARGARÁ DE CLASIFICAR EL CONTENIDO Y MIRAR SI SE AJUSTAN A LOS PATRONES ESTABLECIDOS (AQUÍ ES DONDE PUEDEN SER UTILIZADOS ALGORITMOS COMO LOS VISTOS ANTERIORMENTE). FINALMENTE LAS PÁGINAS QUE HAN CLASIFICADO SON AÑADIDAS A UNA BASE DE DATOS QUE SERVIRÁ DE BASE PARA FUTURAS BÚSQUEDA.

### 3. SELECCIÓN Y DOCUMENTACIÓN DEL SOFTWARE PARA EL MODELO EMERGENTE

PARA INSTALAR RUBY LA MEJOR ALTERNATIVA ES RVM (RUBY VERSION MANAGER):  
[HTTPS://RVM.IO//](https://rvm.io//)

PARA LA REALIZACIÓN DE LA APLICACIÓN UTILIZAREMOS EL LENGUAJE DE PROGRAMACIÓN RUBY, EL FRAMEWORK PARA DESARROLLO DE APLICACIONES WEB RAILS, Y LAS SIGUIENTE LISTA DE LIBRERÍAS:

- NET/HTTP ([HTTP://RUBY-DOC.ORG/STDLIB-1.9.3/LIBDOC/NET/HTTP/RDOC/NET/HTTP.HTML](http://ruby-doc.org/stdlib-1.9.3/libdoc/net/http/rdoc/net/http.html))
- HPRICOT ([HTTPS://GITHUB.COM/HPRICOT/HPRICOT](https://github.com/hpricot/hpricot))

- `JSON(HTTP://DEVELOPER.YAHOO.COM/RUBY/RUBY-JSON.HTML)`
- `OPEN-URI(HTTP://JURETTA.COM/LOG/2006/08/13/RUBY_NET_HTTP_AND_OPEN-URI/)`
- `OMNIAUTH(HTTPS://GITHUB.COM/INTRIDEA/OMNIAUTH)`

TAMBIÉN EXISTE LA POSIBILIDAD DE UTILIZAR LA LIBRERÍA ANEMONE, LA CUAL SE UTILIZA PARA IMPLEMENTAR CRAWLERS CON RUBY. TODO DEPENDE DE LAS DECISIONES DE IMPLEMENTACIÓN QUE TOMEMOS EN EL CAMINO.

A CONTINUACIÓN MOSTRAREMOS LAS PRUEBAS QUE HEMOS REALIZADO CON LAS LIBRERÍAS QUE VAMOS A UTILIZAR:

- EJEMPLO PARA BÚSQUEDA DE DATOS UTILIZANDO LAS LIBRERÍAS LISTADAS ANTERIORMENTE:

**CÓDIGO:**

```

multiplexParser.rb — Ruby
* startParser.rb * multiplexParser.rb * anemone.rb

require 'open-uri'
require 'rubygems'
require 'hpricot'

class MultiplexParser

  @cities_cinecol
  @movies_cinecol
  @cities_provincial

  def initialize site
    if site =~ /cinacolombia/
      @cities_cinecol = []
      @movies_cinecol = []
      url_cinecol = "http://www.cinacolombia.com/Default.aspx"
      save_cities_cinecol url_cinecol
      movies_cinecol url_cinecol
    elsif site =~ /provincial/
      @cities_provincial = []
      url_provincial = "http://provincial.com/index.php"
      save_cities_provincial url_provincial
      movies_provincial url_provincial
    end
  end

  def save_cities_cinecol url
    cinecol = Hpricot(open(url))
    cities = (cinecol/"select.combociudades/option")
    cities.each do |city|
      val = city.inner_html
      unless /Cambiar Ciudad/.match(val) or /----/.match(val)
        @cities_cinecol.push({:id => city.attributes['value'] , :name => val})
      end
    end
    @cities_cinecol.uniq!
    puts "*** Ciudades disponibles en cinacolombia ***"
    for city in @cities_cinecol
      puts "=> ID de película : #{city[:id]} , nombre: #{city[:name]}"
    end
  end
end
```

```

def movies_cinecol url
  cinecol = Hpricot(open(url))
  movies = (cinecol/"select.combopeliculas/option")
  movies.each do |m|
    val = m.inner_html
    @movies_cinecol.push(val) unless /Por Pelicula/.match(val) or /---/.match(val)
  end
  @movies_cinecol.uniq!
  puts "**** Peliculas disponibles en cinecolombia ****"
  for movie in @movies_cinecol
    puts "=> #{movie}"
  end
end

def save_cities_procinal url
  procinal = Hpricot(open(url))
  cities = (procinal/"div.filtrogroup:first ul li a")
  cities.each do |city|
    value = city.inner_html
    @cities_procinal.push(value)
  end
  @cities_procinal.uniq!
  puts "ciudades disponibles en procinal son: #{@cities_procinal}"
end

def movies_procinal url
  procinal = Hpricot(open(url))
  movies = (procinal/"div#cuadriculaPeliculas ul li a")
  #puts movies
end
end

```

## RESULTADOS DE LA EJECUCIÓN

```

MacBook-Pro-dmuneras:PracticaTelemaco danielmunerasanchez$ ruby startParser.rb cinecolombia
*** Ciudades disponibles en cinecolombia ***
=> ID de pelicula : 2 , nombre: ARMENIA
=> ID de pelicula : 3 , nombre: BARRANQUILLA
=> ID de pelicula : 1 , nombre: BOGOT&#193;
=> ID de pelicula : 4 , nombre: BUCARAMANGA
=> ID de pelicula : 5 , nombre: CALI
=> ID de pelicula : 6 , nombre: CARTAGENA
=> ID de pelicula : 24 , nombre: MANIZALES
=> ID de pelicula : 7 , nombre: MEDELL&#205;N
=> ID de pelicula : 8 , nombre: MONTER&#205;A
=> ID de pelicula : 9 , nombre: PEREIRA
**** Peliculas disponibles en cinecolombia ****
=> &#161;PIRATAS! UNA LOCA AVENTURA
=> AMERICAN PIE, EL REENCUENTRO
=> BATALLA NAVAL
=> COMANDO ESPECIAL
=> CUANDO TE ENCUENTRE
=> EL CONSPIRADOR
=> EL ENVIADO
=> EL NI&#209;O DE LA BICICLETA
=> EL PRECIO DE LA CODICIA
=> EL ROMANCE DEL SIGLO
=> EL SUE&#209;O DE IVAN
=> EL VENGADOR FANTASMA 2 3D
=> ESPEJITO, ESPEJITO
=> GORDO, CALVO Y BAJITO
=> HOMBRES DE NEGRO III 3D
=> LORAX - EN BUSCA DE LA TRUFULA PERDIDA
=> LOS VENGADORES
=> MUJERES AL PODER
=> OPERA ON ICE
=> VOTOS DE AMOR
MacBook-Pro-dmuneras:PracticaTelemaco danielmunerasanchez$ █

```

**CÓDIGO EJEMPLO DE ANEMONE:**

```
* startParser.rb * multiplexParser.rb * anemone.rb

require 'anemone'
require 'open-uri'
require 'hpricot'

unless ARGV.length == 1
  STDERR.puts "Usage: #{$0} <url>"
  exit 1
end

$url = ARGV[0]

puts "Pagina a visitar: #{$url}"

def who_is_it site
  if site =~ /cinecolombia/
    puts "Estamos en cine colombia"
  elsif site =~ /cinemark/
    puts "Estamos en cinemark"
  end
end

who_is_it $url

link = 0
Anemone.crawl($url) do |anemone|
  anemone.on_every_page do |page|
    link = link + 1
    puts page.url
  end
end

puts link.to_s
```



```

MacBook-Pro-dmuneras:PracticaTelemaco danielmunerasanchez$ LS
anemone.rb      multiplexParser.rb      startParser.rb
MacBook-Pro-dmuneras:PracticaTelemaco danielmunerasanchez$ ruby anemone.rb http://www.cinecolombia.com/Default.aspx
Pagina a visitar: http://www.cinecolombia.com/Default.aspx
Estamos en cine colombia
http://www.cinecolombia.com/Default.aspx
http://www.cinecolombia.com/contact.aspx
http://www.cinecolombia.com/contactPQRS.aspx
http://www.cinecolombia.com/servicios
http://www.cinecolombia.com/servicios/
http://www.cinecolombia.com/food.aspx
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=814
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=820
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=846
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=812
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=854
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=815&idciudad=1
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=826
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=646
http://www.cinecolombia.com/movie.aspx?idciudad=1&idpelicula=864

```

## Bibliografía

- A. Rungsawang\*, N. A. (7 de Enero de 2004). *Learnable topic-specific web crawler*. Recuperado el 19 de Mayo de 2012, de [http://ac.els-cdn.com/S1084804504000086/1-s2.0-S1084804504000086-main.pdf?\\_tid=0d09e274ee474e734377edf8311a4fce&acdnt=1337445117\\_7dd90c03ab3150abec317e34c96a8a89](http://ac.els-cdn.com/S1084804504000086/1-s2.0-S1084804504000086-main.pdf?_tid=0d09e274ee474e734377edf8311a4fce&acdnt=1337445117_7dd90c03ab3150abec317e34c96a8a89)
- Chang Su, Y. G. (2005). *An Efficient Adaptive Focused Crawler Based on Ontology Learning*. Recuperado el 19 de Mayo de 2012, de <http://csdl.computer.org.ezproxy.eafit.edu.co/dl/proceedings/his/2005/2457/00/24570073.pdf>
- Fatemeh Ahmadi-Abkenari, A. S. (s.f.). *An architecture for a focused trend parallel Web crawler*. Recuperado el 19 de Mayo de 2012, de <http://dl.acm.org/citation.cfm?id=2051686>
- Luján, R. (s.f.). *Crawling the Web*. Recuperado el 19 de Mayo de 2012, de <http://cs.uns.edu.ar/~agm/mineriaweb/downloads/Slides/clase06-slides-rocio.pdf>
- Martínez, F. (s.f.). *InfoSpider*. Recuperado el 19 de Mayo de 2012, de <http://trevinca.ei.uvigo.es/~pcuesta/sm/alumnos2002/infoSpiders.pdf>
- Michael Hersovici, M. J. (s.f.). *The Shark-Search Algorithm*. Recuperado el 2012 de Mayo de 19, de <http://www.cs.cmu.edu/~dpelleg/bin/360.html>

Mukherjea, S. (s.f.). *WTMS: A System for Collecting and Analyzing Topic-Specific Web Information*. Recuperado el 19 de Mayo de 2012, de <http://www9.org/w9cdrom/293/293.html#Muk00>

Yang, S.-Y. (2010). *OntoCrawler: A focused crawler with ontology-supported website models*. Recuperado el 19 de Mayo de 2012, de [http://ac.els-cdn.com/S0957417410000205/1-s2.0-S0957417410000205-main.pdf?\\_tid=9135954b26cd1f0d5a3b8153e0084368&acdnat=1337445082\\_83a0e61583c1bce367b3c71205010057](http://ac.els-cdn.com/S0957417410000205/1-s2.0-S0957417410000205-main.pdf?_tid=9135954b26cd1f0d5a3b8153e0084368&acdnat=1337445082_83a0e61583c1bce367b3c71205010057)