



(19) **United States**

(12) **Patent Application Publication**
Hogue et al.

(10) **Pub. No.: US 2007/0198481 A1**

(43) **Pub. Date: Aug. 23, 2007**

(54) **AUTOMATIC OBJECT REFERENCE IDENTIFICATION AND LINKING IN A BROWSEABLE FACT REPOSITORY**

Publication Classification

(51) **Int. Cl.**
G06F 17/30 (2006.01)

(52) **U.S. Cl.** **707/3**

(76) Inventors: **Andrew W. Hogue**, Ho Ho Kus, NJ (US); **Jonathan T. Betz**, Summit, NJ (US)

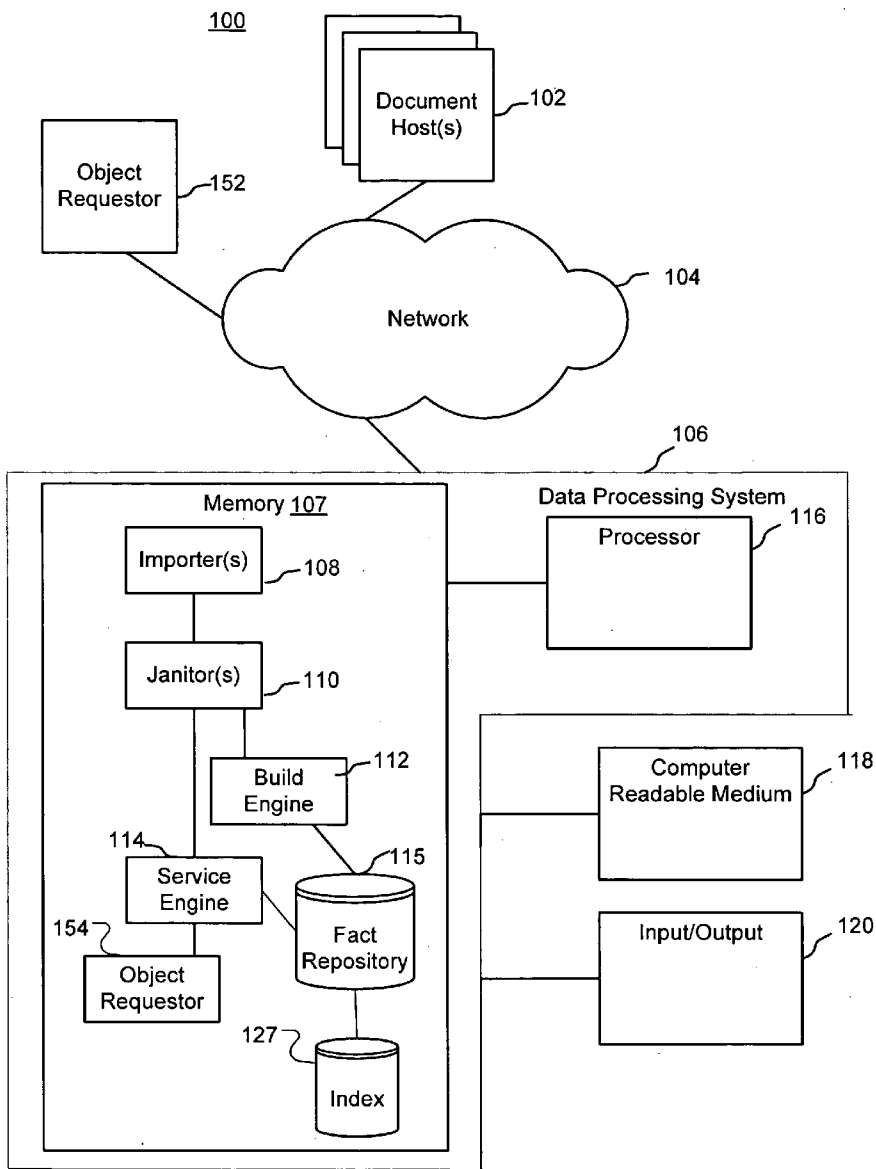
(57) **ABSTRACT**

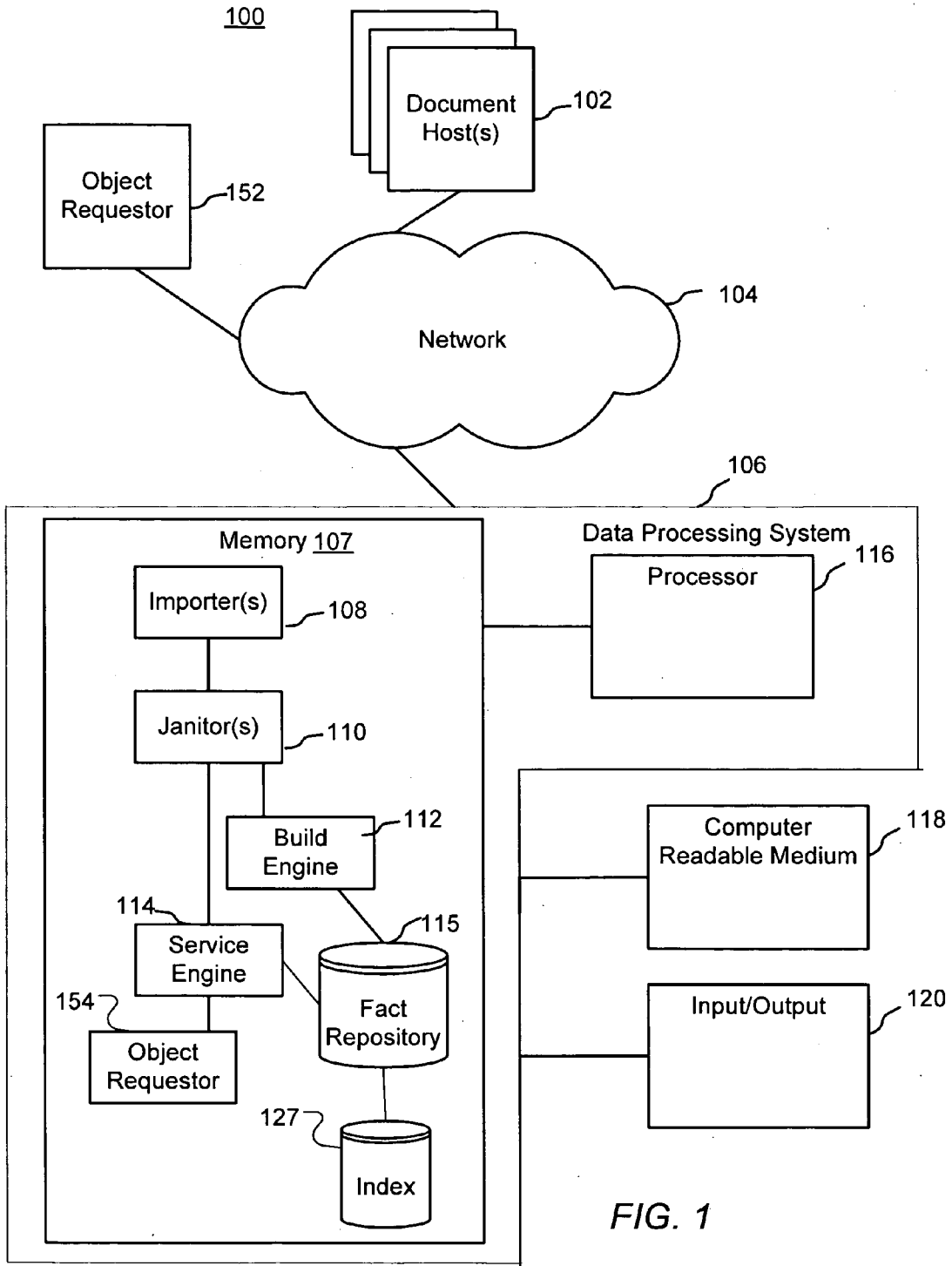
Links between facts associated with objects are automatically created and maintained in a fact repository. Names of objects are automatically identified in the facts, and collected into a list of names. The facts are then processed to identifying such names in the facts. Identified names are used as anchor text for search links. A search link includes a search query for a service engine which search the fact repository for facts associated with objects having the same name.

Correspondence Address:
GOOGLE / FENWICK
SILICON VALLEY CENTER
801 CALIFORNIA ST.
MOUNTAIN VIEW, CA 94041 (US)

(21) Appl. No.: **11/356,837**

(22) Filed: **Feb. 17, 2006**





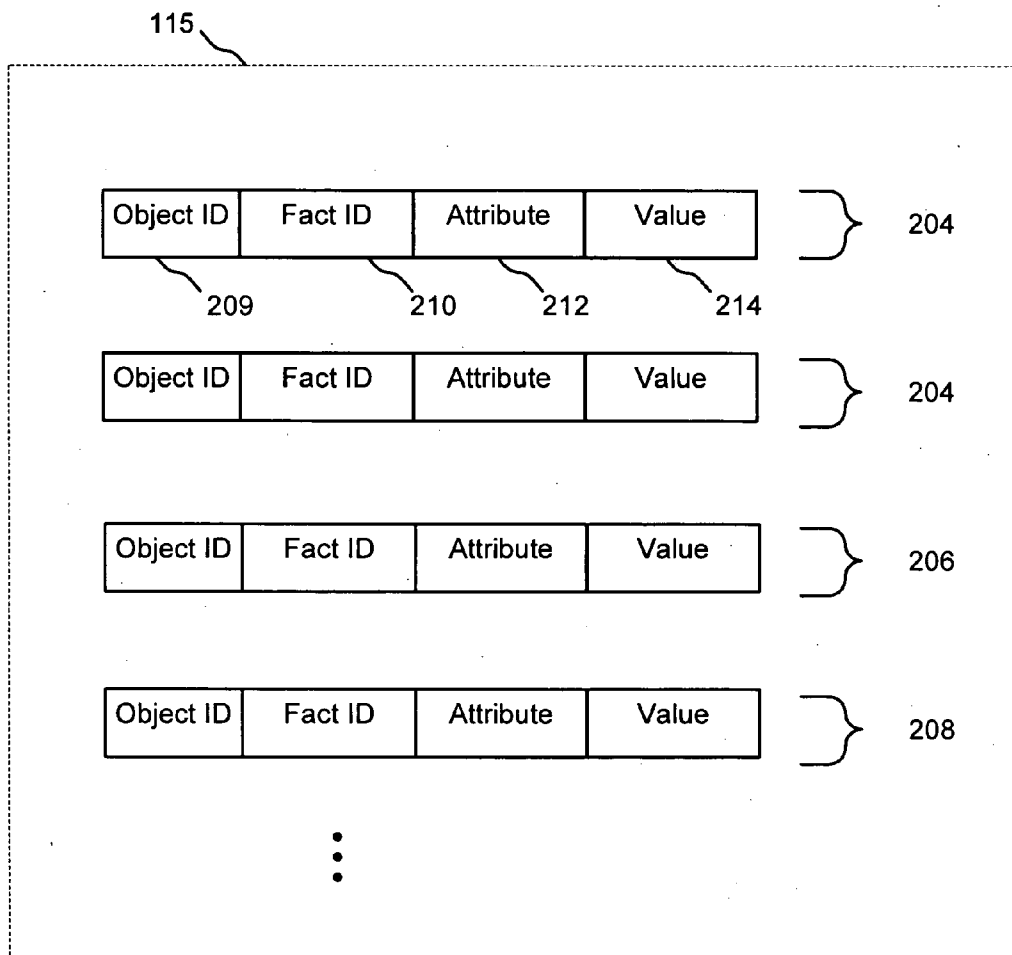
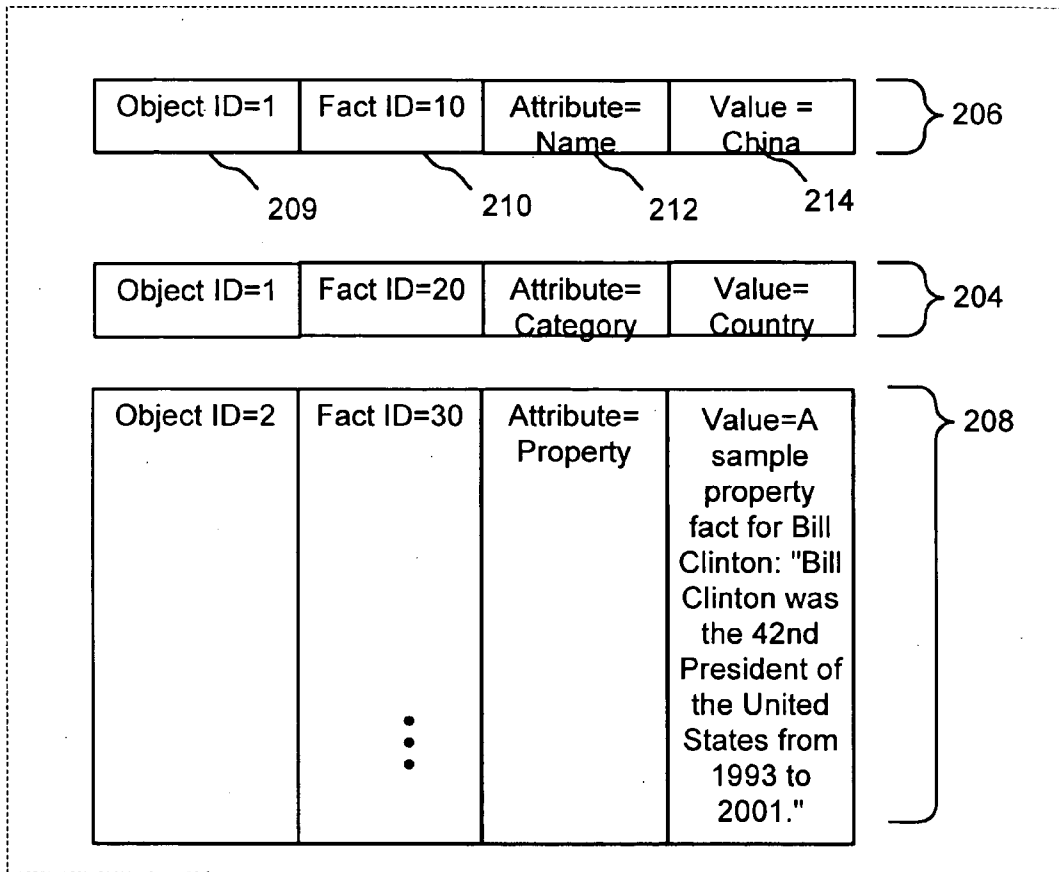


FIG. 2(a)
Example Format of Facts in Repository (each fact is associated with an object ID)



115

FIG. 2(b)

*Example Facts in Repository
(each fact is associated with an
object ID)*

Object ID=1	Fact ID=10
Object ID=1	Fact ID=20
Object ID=1	Fact ID=30
Object ID=2	Fact ID=40

⋮

FIG. 2(c)
*Example Object
Reference Table*

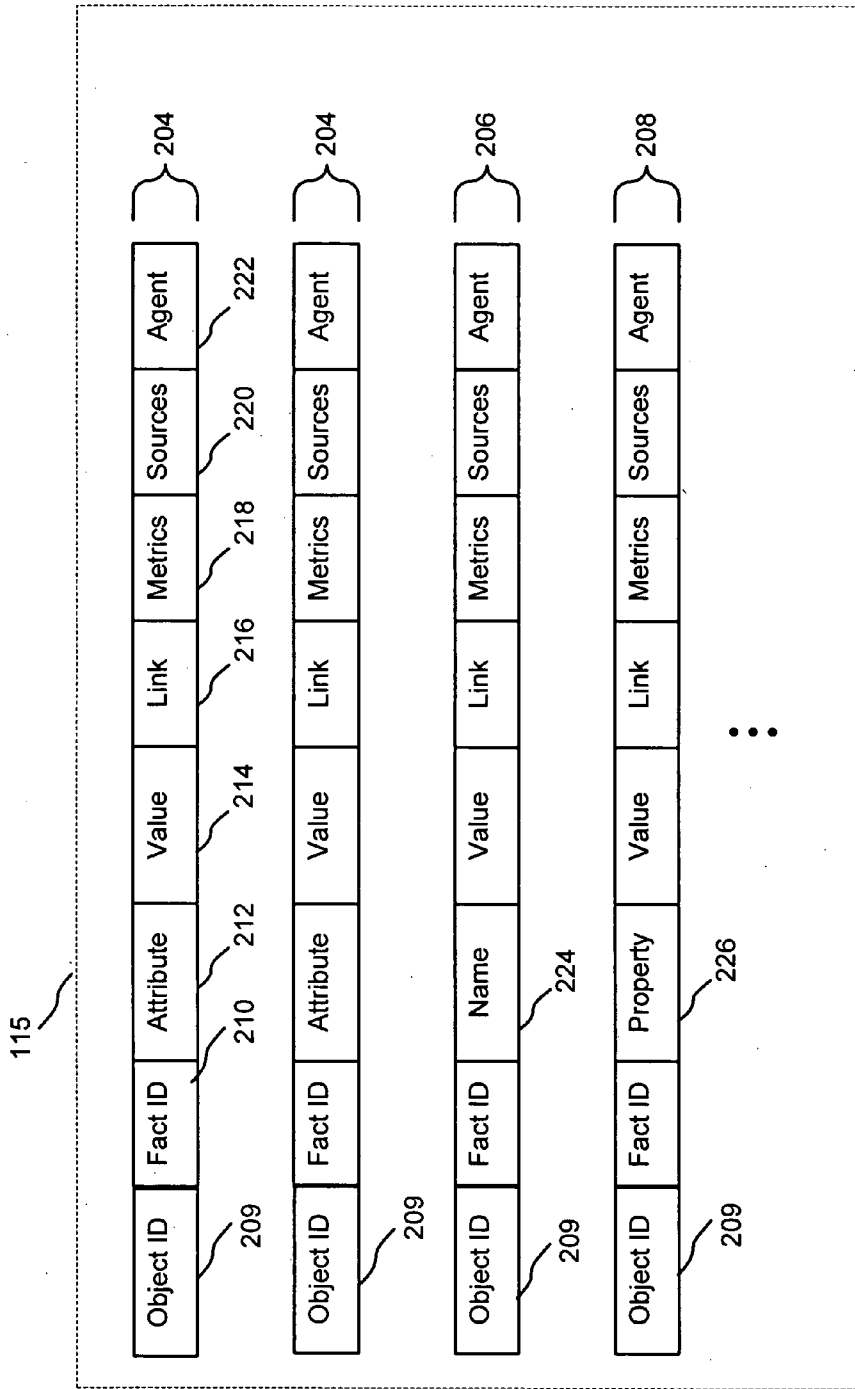


FIG. 2(d)
Example Format of Facts in
Repository (each fact is
associated with an object ID)

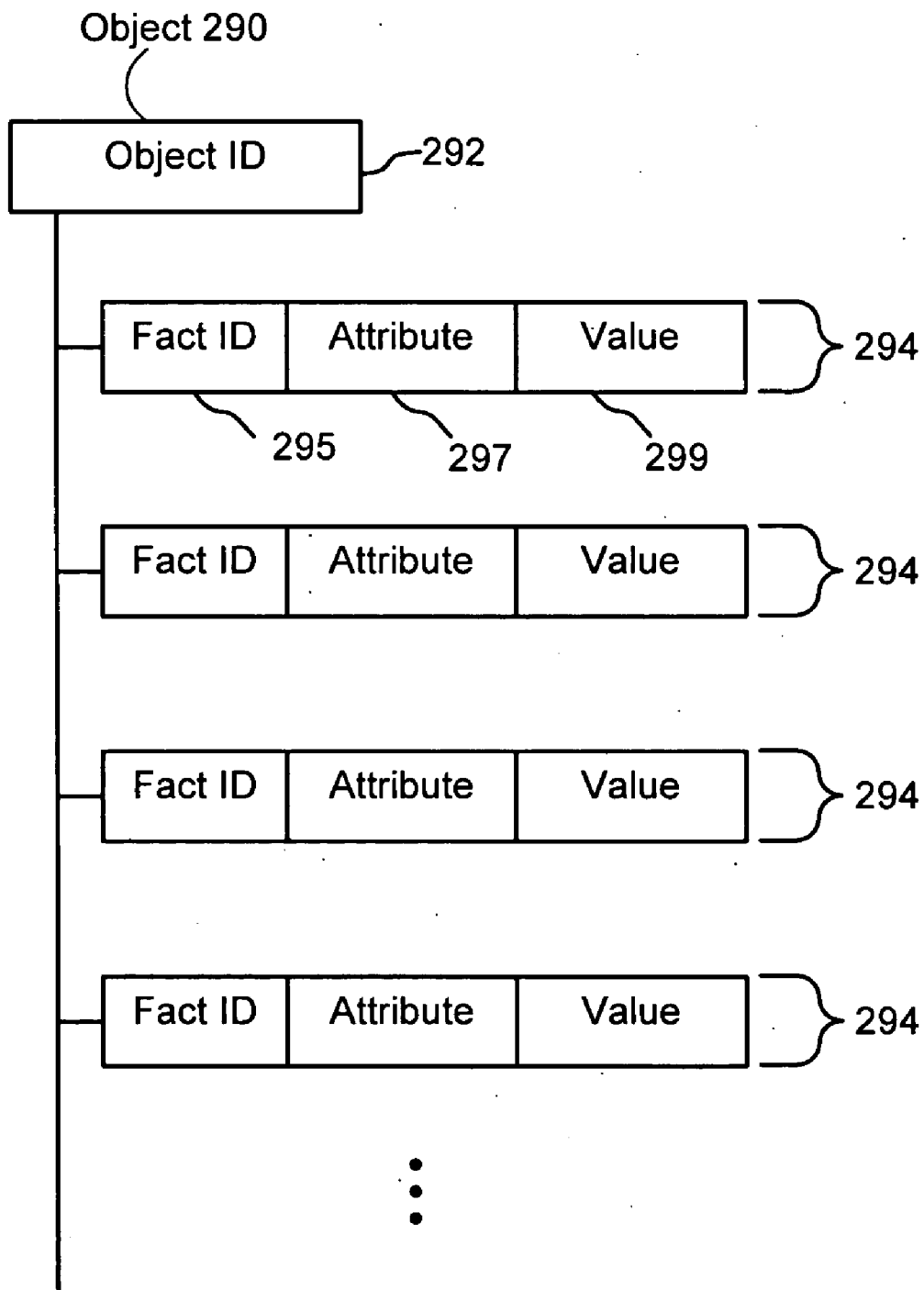


FIG. 2(e)
Example Objects

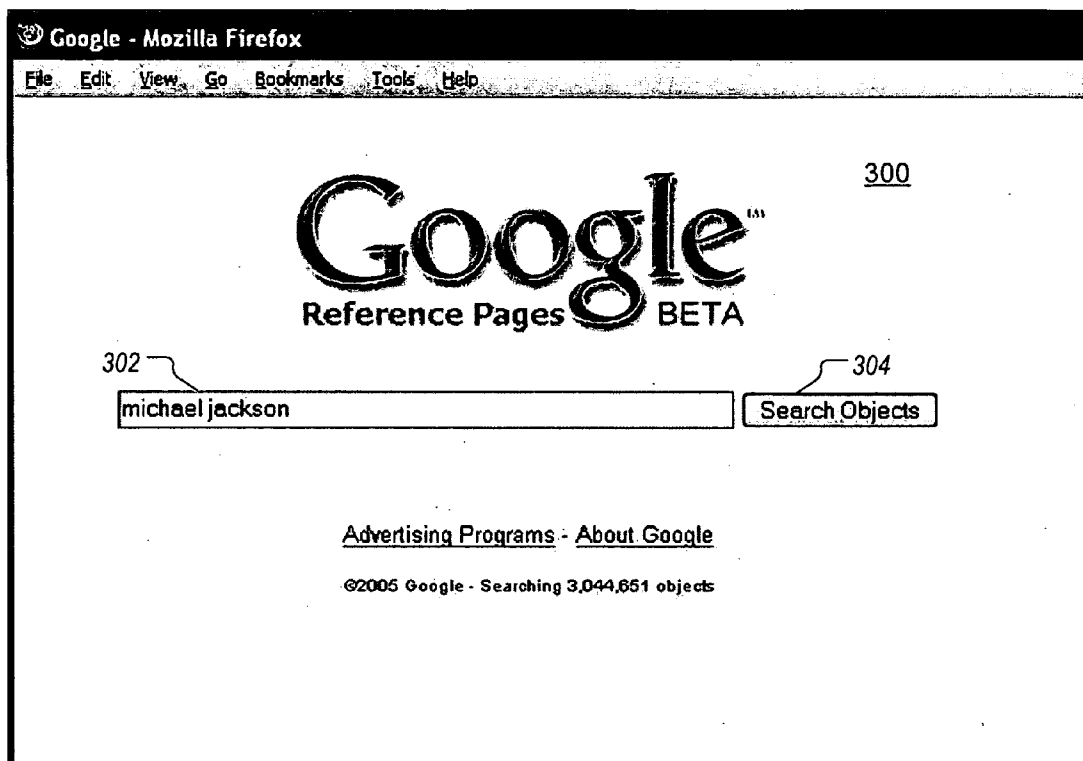


FIG. 3

Google Search: michael jackson - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

Google Reference Pages BETA

michael jackson Search Objects

302 400 304

Reference Pages

402 { Michael Jackson [person] 416 406 410
Date of Birth: 8 November 1970; Place of Birth: Ottawa, Ontario, Canada; Controller, Chief Accounting Officer of: AGENCY.COM Ltd.;...
More from ... >

402 { Michael Jackson [person] 416 406 410
Date of Birth: 29 August 1958; Place of Birth: Gary, Indiana; Full Name: Michael Joseph Jackson;...
More from ... >

402 { Michael Jackson [person] 416 406
Drafted: Drafted by the New York Knickerbockers in the 2nd round (23rd pick) of the 1986 NBA draft; Place of Birth: Fairfax, Virginia;...
More from ... >

402 { Michael Jackson (Anglican bishop) [Anglican bishop] 416 406
Date of Birth: 24 May 1956; Full Name: Michael Jackson;
More from ... >

FIG. 4

Google Search: michael jackson - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

Google Search Objects Search the Web [Advanced Search](#) [Preferences](#)

Reference Pages BETA 406 500

Michael Jackson [person] Results 1 - 1 of 1 for michael jackson (6 ms)

510a Date of Birth 29 August 1958 (According to: [en.wikipedia.org](#)) 512

510b Place of Birth Gary, Indiana (According to: [concerts.ticketmon.com](#))

Full Name Michael Joseph Jackson (According to: [en.wikipedia.org](#))

Full Name Michael Jackson (According to: [simple.wikipedia.org](#))

Birth Name Michael Joseph Jackson II (According to: [www.123musicstars.com](#))

Spouse Lisa Marie Presley (18 May 1994 - 18 January 1996) (divorced) (According to: [teen.com](#))

Spouse Debbie Rowe (15 November 1996 - 8 October 1999) (divorced) 2 children (According to: [www.imdb.com](#))

Place of Birth Gary, Indiana, USA (According to: [www.123musicstars.com](#))

Height 5' 10 (According to: [www.filmovisportal.com](#))

Height 5 FT. 10 in. (According to: [www.athcelebs.com](#))

Height 5 FT. 10 in. (1m78) (According to: [www.filmbug.co.uk](#))

Band The Jackson Five (According to: [www.soundbug.com](#))

Famous for His 1982 Album Thriller (According to: [www.superheroes.com](#))

Height 5' 10" (1.78 M) (According to: [www.imdb.com](#))

Nickname The Gloved One
Wacko Jacko.
Jacko
King Of Pop
MJ (According to: [www.imdb.com](#))

Listed in Category Musicians & Singers (According to: [www.celebritylink.com](#))

Occupation Musician, Actor (According to: [www.filmbug.co.uk](#))

Occupation Musician, Song Writer, Actor (According to: [www.filmovisportal.com](#))

Famous As Pop singer, composer, actor, producer, songwriter (According to: [www.aceahombi.com](#))

Nationality American (According to: [www.filmovisportal.com](#))

Nick Name Jacko; King of Pop (According to: [www.athcelebs.com](#))

Band The Jacksons (According to: [www.filmbug.co.uk](#))

Sex M (According to: [www.filmovisportal.com](#))

Musical Genre Pop, R&B (According to: [www.soundbug.com](#))

Heritage American (According to: [www.superheroes.com](#))

Name Michael Jackson II (According to: [www.filmovisportal.com](#))

Sign Virgo (According to: [www.filmbug.co.uk](#))

Date of Birth 29 August 1958 (47) (According to: [www.filmbug.co.uk](#))

Born 29 August 1958 (Virgo) (According to: [www.soundbug.com](#))

514 Nationality US - United States of America (According to: [www.athcelebs.com](#))

510c Born As Michael Joseph Jackson (According to: [www.athcelebs.com](#))

Wrote Thriller (1983) (M) (According to: [imagine-contact.com](#))

Wrote Michael Jackson: HISTORY on Film - Volume II (1997) (M) (According to: [www.imdb.com](#))

Wrote Moonwalker (1988) (According to: [musicstars.mymmedia.com](#))

Done

FIG. 5

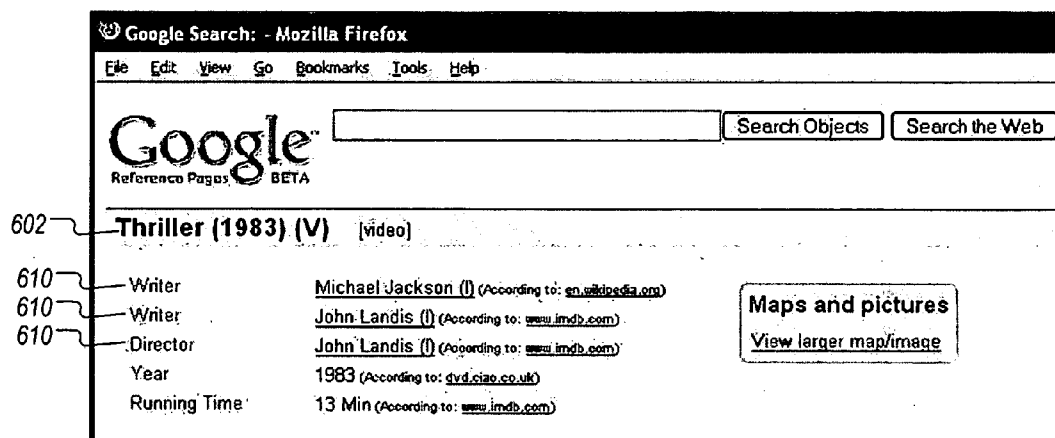


FIG. 6

AUTOMATIC OBJECT REFERENCE IDENTIFICATION AND LINKING IN A BROWSEABLE FACT REPOSITORY

CROSS REFERENCE TO RELATED APPLICATIONS

[0001] This application is related to the following U.S. Applications all of which are incorporated by reference herein:

[0002] U.S. application Ser. No. _____, entitled "Support for Object Search", filed concurrently herewith, by Alex Kehlenbeck, Andrew W. Hogue, Jonathan T. Betz, Attorney Docket No. 24207-10945;

[0003] U.S. application Ser. No. _____, entitled "Data Object Visualization", filed on Jan. 27, 2006, by Andrew W. Hogue, David Vespe, Alex Kehlenbeck, Mike Gordon, Jeffrey C. Reynar, David Alpert, Attorney Docket No. 24207-10946;

[0004] U.S. application Ser. No. _____, entitled "Data Object Visualization Using Maps", filed on Jan. 27, 2006, by Andrew W. Hogue, David Vespe, Alex Kehlenbeck, Mike Gordon, Jeffrey C. Reynar, David Alpert, Attorney Docket No. 24207-10947;

[0005] U.S. application Ser. No. _____, entitled "Query Language", filed concurrently herewith, by Andrew W. Hogue, Doug Rhode, Attorney Docket No. 24207-10948;

[0006] U.S. application Ser. No. _____, entitled "ID Persistence Through Normalization", filed concurrently herewith, by Jonathan T. Betz, Andrew W. Hogue, Attorney Docket No. 24207-10950;

[0007] U.S. application Ser. No. _____, entitled "Browseable Fact Repository", filed concurrently herewith by Andrew W. Hogue, Jonathan T. Betz, Attorney Docket No. 24207-10949;

[0008] U.S. application Ser. No. _____, entitled "Annotation Framework", filed concurrently herewith, by Tom Richford, Jonathan T. Betz, Attorney Docket No. 24207-10951;

[0009] U.S. application Ser. No. _____, entitled "Object Categorization for Information Extraction", filed on Jan. 27, 2006, by Jonathan T. Betz, Attorney Docket No. 24207-10952;

[0010] U.S. application Ser. No. _____, entitled "Modular Architecture for Entity Normalization", filed concurrently herewith, by Jonathan T. Betz, Farhan Shamsi, Attorney Docket No. 24207-10953;

[0011] U.S. application Ser. No. _____, entitled "Attribute Entropy as a Signal in Object Normalization", filed concurrently herewith by, Jonathan T. Betz, Vivek Menezes, Attorney Docket No. 24207-10954;

[0012] U.S. application Ser. No. _____, entitled "Designating Data Objects for Analysis", filed on Jan. 27, 2006, by Andrew W. Hogue, David Vespe, Alex Kehlenbeck, Mike Gordon, Jeffrey C. Reynar, David Alpert, Attorney Docket No. 24207-11126;

[0013] U.S. application Ser. No. _____, entitled "Data Object Visualization Using Graphs", filed on Jan. 27, 2006,

by Andrew W. Hogue, David Vespe, Alex Kehlenbeck, Mike Gordon, Jeffrey C. Reynar, David Alpert, Attorney Docket No. 24207-11125.

FIELD OF THE INVENTION

[0014] The present invention relates to a knowledge and fact databases, and in particularly to user interfaces, methods, and systems for browsing such database.

BACKGROUND OF THE INVENTION

[0015] Knowledge bases are collections of facts and data organized into systematic arrangements of information. Knowledge bases often include search tools for searching and browsing facts. Online knowledge bases have become increasingly prevalent on the Internet, and examples include WordNet, Wikipedia, Webopedia, and similar online encyclopedias, dictionaries, and document collections.

[0016] In a typical online knowledge base such as Wikipedia or Wordnet, it is common for the content of one entry to include terms or words that are the names or titles or other entries in the database. Thus, in the Wikipedia online encyclopedia, the entry for "American Revolution" includes many terms, such as "George Washington" or "Declaration of Independence." Given that these databases use hypertext, it is conventional that a phrase such as "George Washington" is constructed as the anchor text for a hyperlink from the "American Revolution" entry to the "George Washington" entry. More specifically, these types of hyperlinks are fixed in that they reference a specific page or entry in the database. For example, in the Wikipedia entry for "American Revolution", the hyperlink for the phrase "George Washington" is the following URL (uniform resource locator):

[0017] [http://en.wikipedia.org/wiki/George Washington](http://en.wikipedia.org/wiki/George_Washington)

[0018] where the last portion of the URL, "George_Washington" identifies a specific, pre-existing document in the Wikipedia document collection.

[0019] This approach to linking of phrases between entries in an online knowledge base works only where the referenced entry (e.g., here the entry for "George Washington") exists at the time the referencing entry (e.g., here the entry for "American Revolution") is written, so that the latter's author can manually include the appropriate link to the referenced entry.

[0020] A problem arises then in an online knowledge base which is under going frequent changes, including the addition of new entries, and changes in existing entries. An existing entry may include terms or phrases for which there were no corresponding entries at the time the existing entry was written. For example, a knowledge base may include an entry on quantum mechanics, and a reference to string theory, but at the time quantum mechanics entry is written there is no other entry that describes string theory. However, at a late date, new entries may have been added which could be properly referenced for the terms of the existing entry. Thus, a new, subsequent entry in the database with the title "String Theory" may be created. Since these terms were not linked at the time the existing entry was authored, a user reading the entry on quantum mechanics would not know there is a corresponding entry on string theory.

[0021] Another problem with fixed (or "hard") entries is that they make it more difficult to update the database or

change its file structure, since the pathname name of the referenced article cannot be changed without causing the hyperlink to malfunction.

SUMMARY OF THE INVENTION

[0022] The present invention provides a methodology and system for automatically creating and maintaining links between facts in a fact repository. The fact repository includes a large collection of facts, each of which is associated with an object, such as a person, place, book, movie, country, or any other entity of interest. Each fact comprises an attribute, which is descriptive of the type of fact (e.g., "name," or "population"), and a value for that attribute (e.g., "George Washington", or "1,397,264,580"). A value can also contain any amount of text—from a single term or phrase to many paragraphs or pages—such as appropriate to describe the attribute. Each object will have a name fact that is the name of the object. The value of a value can thus include one or more phrases that are themselves the names of other facts.

[0023] In one embodiment, a two stage process is used to automatically construct and maintain links between facts in the repository. In a first stage, a set of name facts is collected to form a list of the name of objects in the repository. In a next stage, the value portion of some set of facts is processed to identify whether they include one or more names on the list of object names. (This second set may be the same set or a different set than that from which the object names were collected.) Where an object name is identified in the value of a fact, a search link is constructed, using the object name as anchor text, but using a search query for the link contents. The search query is a query for any objects which have as a name the object name found in the value of the fact. Constructing the search link in this manner ensures that all facts that include the name of other facts will be properly linked to such other facts, even if they are created subsequently.

[0024] In one embodiment, the identification of object names in the fact values uses a phrase identification algorithm to first identify phrases in text of the fact values. Next, each of the identified phrases is compared with the list of object names to determine if it appears on the list. If an identified phrase is present in the list of object names, then the phrase is used as the anchor text for the search link.

[0025] The foregoing process is preferably repeated on periodic basis, for example daily, weekly, or any other time period. This ensures that new facts entered in the interval between repetitions are processed both to include the new object names on the object name list (so that existing facts can now properly reference such new object names), and to ensure that objects named in the value of such new facts are linked to previously existing facts.

[0026] The present invention further has embodiments in computer program products, in computer systems, and computer user interfaces, which various perform or cooperate in the operation or use of the foregoing method (or its alternative embodiments and features).

BRIEF DESCRIPTION OF THE DRAWINGS

[0027] FIG. 1 shows a system architecture, in accordance with a preferred embodiment of the invention.

[0028] FIGS. 2(a)-2(d) are block diagrams illustrating a data structure for facts within a repository of FIG. 1 in accordance with preferred embodiments of the invention.

[0029] FIG. 2(e) is a block diagram illustrating an alternate data structure for facts and objects in accordance with preferred embodiments of the invention.

[0030] FIG. 3 illustrates a landing page for initiating a search query of a fact repository.

[0031] FIG. 4 illustrates a search results page for a search on the fact repository.

[0032] FIG. 5 illustrates an object detail page.

[0033] FIG. 6 illustrates another object detail page obtained in response to selection of the object reference link in FIG. 5.

[0034] The figures depict a preferred embodiment of the present invention for purposes of illustration only. One skilled in the art will readily recognize from the following discussion that alternative embodiments of the structures and methods illustrated herein may be employed without departing from the principles of the invention described herein.

DETAILED DESCRIPTION

System Overview

[0035] FIG. 1 shows a system architecture 100 adapted to support one embodiment of the invention. FIG. 1 shows components used to add facts into, and retrieve facts from a repository 115. The system architecture 100 includes a network 104, through which any number of document hosts 102 communicate with a data processing system 106, along with any number of object requesters 152, 154.

[0036] Document hosts 102 store documents and provide access to documents. A document is comprised of any machine-readable data including any combination of text, graphics, multimedia content, etc. A document may be encoded in a markup language, such as Hypertext Markup Language (HTML), i.e., a web page, in a interpreted language (e.g., JavaScript) or in any other computer readable or executable format. A document can include one or more hyperlinks to other documents. A typical document will include one or more facts within its content. A document stored in a document host 102 may be located and/or identified by a Uniform Resource Locator (URL), or Web address, or any other appropriate form of identification and/or location. A document host 102 is implemented by a computer system, and typically includes a server adapted to communicate over the network 104 via networking protocols (e.g., TCP/IP), as well as application and presentation protocols (e.g., HTTP, HTML, SOAP, D-HTML, Java). The documents stored by a host 102 are typically held in a file directory, a database, or other data repository. A host 102 can be implemented in any computing device (e.g., from a PDA or personal computer, a workstation, mini-computer, or mainframe, to a cluster or grid of computers), as well as in any processor architecture or operating system.

[0037] FIG. 1 shows components used to manage facts in a fact repository 115. Data processing system 106 includes one or more importers 108, one or more janitors 110, a build engine 112, a service engine 114, and a fact repository 115

(also called simply a “repository”). Each of the foregoing are implemented, in one embodiment, as software modules (or programs) executed by processor **116**. Importers **108** operate to process documents received from the document hosts, read the data content of documents, and extract facts (as operationally and programmatically defined within the data processing system **106**) from such documents. The importers **108** also determine the subject or subjects with which the facts are associated, and extract such facts into individual items of data, for storage in the fact repository **115**. In one embodiment, there are different types of importers **108** for different types of documents, for example, dependent on the format or-document type.

[0038] Janitors **110** operate to process facts extracted by importer **108**. This processing can include but is not limited to, data cleansing, object merging, and fact induction. In one embodiment, there are a number of different janitors **110** that perform different types of data management operations on the facts. For example, one janitor **110** may traverse some set of facts in the repository **115** to find duplicate facts (that is, facts that convey the same factual information) and merge them. Another janitor **110** may also normalize facts into standard formats. Another janitor **110** may also remove unwanted facts from repository **115**, such as facts related to pornographic content. Other types of janitors **110** may be implemented, depending on the types of data management functions desired, such as translation, compression, spelling or grammar correction, and the like.

[0039] Various janitors **110** act on facts to normalize attribute names, and values and delete duplicate and near-duplicate facts so an object does not have redundant information. For example, we might find on one page that Britney Spears’ birthday is “Dec. 2, 1981” while on another page that her date of birth is “Dec. 2, 1981.” Birthday and Date of Birth might both be rewritten as Birthdate by one janitor and then another janitor might notice that Dec. 2, 1981 and Dec. 2, 1981 are different forms of the same date. It would choose the preferred form, remove the other fact and combine the source lists for the two facts. As a result when you look at the source pages for this fact, on some you’ll find an exact match of the fact and on others text that is considered to be synonymous with the fact.

[0040] Build engine **112** builds and manages the repository **115**. Service engine **114** is an interface for querying the repository **115**. Service engine **114**’s main function is to process queries, score matching objects, and return them to the caller but it is also used by janitor **110**.

[0041] Repository **115** stores factual information extracted from a plurality of documents that are located on document hosts **102**. A document from which a particular fact may be extracted is a source document (or “source”) of that particular fact. In other words, a source of a fact includes that fact (or a synonymous fact) within its contents.

[0042] Repository **115** contains one or more facts. In one embodiment, each fact is associated with exactly one object. One implementation for this association includes in each fact an object ID that uniquely identifies the object of the association. In this manner, any number of facts may be associated with an individual object, by including the object ID for that object in the facts. In one embodiment, objects themselves are not physically stored in the repository **115**, but rather are defined by the set or group of facts with the

same associated object ID, as described below. Further details about facts in repository **115** are described below, in relation to FIGS. **2(a)**-**2(d)**.

[0043] It should be appreciated that in practice at least some of the components of the data processing system **106** will be distributed over multiple computers, communicating over a network. For example, repository **115** may be deployed over multiple servers. As another example, the janitors **110** may be located on any number of different computers. For convenience of explanation, however, the components of the data processing system **106** are discussed as though they were implemented on a single computer.

[0044] In another embodiment, some or all of document hosts **102** are located on data processing system **106** instead of being coupled to data processing system **106** by a network. For example, importer **108** may import facts from a database that is a part of or associated with data processing system **106**.

[0045] FIG. **1** also includes components to access repository **115** on behalf of one or more object requesters **152**, **154**. Object requesters are entities that request objects from repository **115**. Object requesters **152**, **154** may be understood as clients of the system **106**, and can be implemented in any computer device or architecture. As shown in FIG. **1**, a first object requester **152** is located remotely from system **106**, while a second object requester **154** is located in data processing system **106**. For example, in a computer system hosting a blog, the blog may include a reference to an object whose facts are in repository **115**. An object requester **152**, such as a browser displaying the blog will access data processing system **106** so that the information of the facts associated with the object can be displayed as part of the blog web page. As a second example, janitor **110** or other entity considered to be part of data processing system **106** can function as object requester **154**, requesting the facts of objects from repository **115**.

[0046] FIG. **1** shows that data processing system **106** includes a memory **107** and one or more processors **116**. Memory **107** includes importers **108**, janitors **110**, build engine **112**, service engine **114**, and requester **154**, each of which are preferably implemented as instructions stored in memory **107** and executable by processor **116**. Memory **107** also includes repository **115**. Repository **115** can be stored in a memory of one or more computer systems or in a type of memory such as a disk. FIG. **1** also includes a computer readable medium **118** containing, for example, at least one of importers **108**, janitors **110**, build engine **112**, service engine **114**, requester **154**, and at least some portions of repository **115**. FIG. **1** also includes one or more input/output devices **120** that allow data to be input and output to and from data processing system **106**. It will be understood that data processing system **106** preferably also includes standard software components such as operating systems and the like and further preferably includes standard hardware components not shown in the figure for clarity of example.

[0047] FIG. **2(a)** shows an example format of a data structure for facts within repository **115**, according to some embodiments of the invention. As described above, the repository **115** includes facts **204**. Each fact **204** includes a unique identifier for that fact, such as a fact ID **210**. Each fact **204** includes at least an attribute **212** and a value **214**.

For example, a fact associated with an object representing George Washington may include an attribute of “date of birth” and a value of “Feb. 22, 1732.” In one embodiment, all facts are stored as alphanumeric characters since they are extracted from web pages. In another embodiment, facts also can store binary data values. Other embodiments, however, may store fact values as mixed types, or in encoded formats.

[0048] As described above, each fact is associated with an object ID **209** that identifies the object that the fact describes. Thus, each fact that is associated with a same entity (such as George Washington), will have the same object ID **209**. In one embodiment, objects are not stored as separate data entities in memory. In this embodiment, the facts associated with an object contain the same object ID, but no physical object exists. In another embodiment, objects are stored as data entities in memory, and include references (for example, pointers or IDs) to the facts associated with the object. The logical data structure of a fact can take various forms; in general, a fact is represented by a tuple that includes a fact ID, an attribute, a value, and an object ID. The storage implementation of a fact can be in any underlying physical data structure.

[0049] FIG. 2(b) shows an example of facts having respective fact IDs of **10**, **20**, and **30** in repository **115**. Facts **10** and **20** are associated with an object identified by object ID “**1**.” Fact **10** has an attribute of “Name” and a value of “China.” Fact **20** has an attribute of “Category” and a value of “Country.” Thus, the object identified by object ID “**1**” has a name fact **205** with a value of “China” and a category fact **206** with a value of “Country.” Fact **30208** has an attribute of “Property” and a value of “Bill Clinton was the 42nd President of the United States from 1993 to 2001.” Thus, the object identified by object ID “**2**” has a property fact with a fact ID of **30** and a value of “Bill Clinton was the 42nd President of the United States from 1993 to 2001.” In the illustrated embodiment, each fact has one attribute and one value. The number of facts associated with an object is not limited; thus while only two facts are shown for the “China” object, in practice there may be dozens, even hundreds of facts associated with a given object. Also, the value fields of a fact need not be limited in size or content. For example, a fact about the economy of “China” with an attribute of “Economy” would have a value including several paragraphs of text, numbers, perhaps even tables of figures. This content can be formatted, for example, in a markup language. For example, a fact having an attribute “original html” might have a value of the original html text taken from the source web page.

[0050] Also, while the illustration of FIG. 2(b) shows the explicit coding of object ID, fact ID, attribute, and value, in practice the content of the fact can be implicitly coded as well (e.g., the first field being the object ID, the second field being the fact ID, the third field being the attribute, and the fourth field being the value). Other fields include but are not limited to: the language used to state the fact (English, etc.), how important the fact is, the source of the fact, a confidence value for the fact, and so on.

[0051] FIG. 2(c) shows an example object reference table **210** that is used in some embodiments. Not all embodiments include an object reference table. The object reference table **210** functions to efficiently maintain the associations between object IDs and fact IDs. In the absence of an object

reference table **210**, it is also possible to find all facts for a given object ID by querying the repository to find all facts with a particular object ID. While FIGS. 2(b) and 2(c) illustrate the object reference table **210** with explicit coding of object and fact IDs, the table also may contain just the ID values themselves in column or pair-wise arrangements.

[0052] FIG. 2(d) shows an example of a data structure for facts within repository **115**, according to some embodiments of the invention showing an extended format of facts. In this example, the fields include an object reference link **216** to another object. The object reference link **216** can be an object ID of another object in the repository **115**, or a reference to the location (e.g., table row) for the object in the object reference table **210**. The object reference link **216** allows facts to have as values other objects. For example, for an object “United States,” there may be a fact with the attribute of “president” and the value of “George W. Bush,” with “George W. Bush” being an object having its own facts in repository **115**. In some embodiments, the value field **214** stores the name of the linked object and the link **216** stores the object identifier of the linked object. Thus, this “president” fact would include the value **214** of “George W. Bush”, and object reference link **216** that contains the object ID for the “George W. Bush” object. In some other embodiments, facts **204** do not include a link field **216** because the value **214** of a fact **204** may store a link to another object.

[0053] Each fact **204** also may include one or more metrics **218**. A metric provides an indication of the some quality of the fact. In some embodiments, the metrics include a confidence level and an importance level. The confidence level indicates the likelihood that the fact is correct. The importance level indicates the relevance of the fact to the object, compared to other facts for the same object. The importance level may optionally be viewed as a measure of how vital a fact is to an understanding of the entity or concept represented by the object.

[0054] Each fact **204** includes a list of one or more sources **220** that include the fact and from which the fact was extracted. Each source may be identified by a Uniform Resource Locator (URL), or Web address, or any other appropriate form of identification and/or location, such as a unique document identifier.

[0055] The facts illustrated in FIG. 2(d) include an agent field **222** that identifies the importer **108** that extracted the fact. For example, the importer **108** may be a specialized importer that extracts facts from a specific source (e.g., the pages of a particular web site, or family of web sites) or type of source (e.g., web pages that present factual information in tabular form), or an importer **108** that extracts facts from free text in documents throughout the Web, and so forth.

[0056] Some embodiments include one or more specialized facts, such as a name fact **207** and a property fact **208**. A name fact **207** is a fact that conveys a name for the entity or concept represented by the object ID. A name fact **207** includes an attribute **224** of “name” and a value, which is the name of the object. For example, for an object representing the country Spain, a name fact would have the value “Spain.” A name fact **207**, being a special instance of a general fact **204**, includes the same fields as any other fact **204**; it has an attribute, a value, a fact ID, metrics, sources, etc. The attribute **224** of a name fact **207** indicates that the

fact is a name fact, and the value is the actual name. The name may be a string of characters. An object ID may have one or more associated name facts, as many entities or concepts can have more than one name. For example, an object ID representing Spain may have associated name facts conveying the country's common name "Spain" and the official name "Kingdom of Spain." As another example, an object ID representing the U.S. Patent and Trademark Office may have associated name facts conveying the agency's acronyms "PTO" and "USPTO" as well as the official name "United States Patent and Trademark Office." If an object does have more than one associated name fact, one of the name facts may be designated as a primary name and other name facts may be designated as secondary names, either implicitly or explicitly.

[0057] A property fact 208 is a fact that conveys a statement about the entity or concept represented by the object ID. Property facts are generally used for summary information about an object. A property fact 208, being a special instance of a general fact 204, also includes the same parameters (such as attribute, value, fact ID, etc.) as other facts 204. The attribute field 226 of a property fact 208 indicates that the fact is a property fact (e.g., attribute is "property") and the value is a string of text that conveys the statement of interest. For example, for the object ID representing Bill Clinton, the value of a property fact may be the text string "Bill Clinton was the 42nd President of the United States from 1993 to 2001." Some object IDs may have one or more associated property facts while other objects may have no associated property facts. It should be appreciated that the data structures shown in FIGS. 2(a)-2(d) and described above are merely exemplary. The data structure of the repository 115 may take on other forms. Other fields may be included in facts and some of the fields described above may be omitted. Additionally, each object ID may have additional special facts aside from name facts and property facts, such as facts conveying a type or category (for example, person, place, movie, actor, organization, etc.) for categorizing the entity or concept represented by the object ID. In some embodiments, an object's name(s) and/or properties may be represented by special records that have a different format than the general facts records 204.

[0058] As described previously, a collection of facts is associated with an object ID of an object. An object may become a null or empty object when facts are disassociated from the object. A null object can arise in a number of different ways. One type of null object is an object that has had all of its facts (including name facts) removed, leaving no facts associated with its object ID. Another type of null object is an object that has all of its associated facts other than name facts removed, leaving only its name fact(s). Alternatively, the object may be a null object only if all of its associated name facts are removed. A null object represents an entity or concept for which the data processing system 106 has no factual information and, as far as the data processing system 106 is concerned, does not exist. In some embodiments, facts of a null object may be left in the repository 115, but have their object ID values cleared (or have their importance to a negative value). However, the facts of the null object are treated as if they were removed from the repository 115. In some other embodiments, facts of null objects are physically removed from repository 115.

[0059] FIG. 2(e) is a block diagram illustrating an alternate data structure 290 for facts and objects in accordance with preferred embodiments of the invention. In this data structure, an object 290 contains an object ID 292 and references or points to facts 294. Each fact includes a fact ID 295, an attribute 297, and a value 299. In this embodiment, an object 290 actually exists in memory 107.

[0060] Referring again to FIG. 1, the content of the facts in the repository 115 are also indexed in index 127. The index 127 maintains a term index, which maps terms to {object, fact, field, token} tuples, where "field" is, e.g., attribute or value. The service engine 114 is adapted to receive keyword queries from clients such as object requestors, and communicates with the index 127 to retrieve the facts that are relevant to user's search query. For a generic query containing one or more terms, the service engine 114 assumes the scope is at the object level. Thus, any object with one or more of the query terms somewhere (not necessarily on the same fact) will match the query for purposes of being ranked in the search results.

[0061] In one embodiment the ranking (score) of an object is a linear combination of relevance scores for each of the facts. The relevance score for each fact is based on whether the fact includes one or more query terms (a hit) in either the attribute or value portion of the fact. Each hit is scored based on the frequency of the term that is hit, with more common terms getting lower scores, and rarer terms getting higher scores (e.g., using a TF-IDF based term weighting model). The fact score is then adjusted based on additional factors. These factors include the appearance of consecutive query terms in a fact, the appearance of consecutive query terms in a fact in the order in which they appear in the query, the appearance of an exact match for the entire query, the appearance of the query terms in the name fact (or other designated fact, e.g., property or category), and the percentage of facts of the object containing at least one query term. Each fact's score is also adjusted by its associated confidence measure and by its importance measure. Since each fact is independently scored, the facts most relevant and important to any individual query can be determined, and selected. In one embodiment, a selected number (e.g., 5) of the top scoring facts is selected for display in response to query.

[0062] The service engine 114 is also adapted to handle structured queries, using query operators that restrict the scope of a term match. A fact restriction operator, comprising brackets enclosing terms, e.g., "[terms]", restricts the term(s) to matching in a single fact. Field restriction operators attribute{} and value{} operators restrict to a single field. Thus, the query [attribute{name} value{"george washington"}] restricts the scope of the match to the terms "george washington" appearing in the name fact.

User Interface for Browsing Fact Repository

[0063] Referring now to FIGS. 3-6, there are shown various screens from a user interface for browsing the fact repository 115 in accordance with one embodiment of present invention. In FIG. 3 there is shown a simple landing or home page for initiating a query of the fact repository 115. The page 300 a search query field 302 and a search objects button 304. The user enters any number of search terms into the search field 302, such as the illustrated terms "michael jackson." The terms can be any terms whatsoever, as they

may appear in any fact. The user selects the search objects button **304** to provide the search terms to the service engine **114**.

[0064] FIG. 4 illustrates the search result page **400** of a search, here for the search query “china.” The results page **400** includes a list of ranked search results **402**, each search result **402** comprising a name link **416** to an object, the anchor text of the link **416** being the name of the object (the name link **416** resolves to an object detail page, as further described below). The results **402** are ranked according to their relevance to the search query. Each search result **402** (which for the purpose of this discussion can also be referred to as an object) is displayed with a label **406** indicating the category of the object (e.g., country, company, person, car, etc.).

[0065] Displayed in conjunction with each search result **402** is a list of one or more facts **410** associated with the object. In this example, the results **402** are for a number of different people named “Michael Jackson”. The list of facts **410** is ordered according to the relevance scores of the facts, so that the fact most relevance to the query are displayed first.

[0066] Selection of the name link **416** for an object search result **402** produces an object detail page **500** for the object, such as illustrated in FIG. 5. The object detail page **500** contains a list of the facts **510** associated with the object, each fact **510** including its attribute (e.g., “name,” “category,” “population,” etc.) and its value. The facts **510** may be listed on a single detailed result page **500**, or on several linked pages. The object detail page **500** can be presented as a table with fixed size columns, or a list. Facts **510** can be textual or graphical. Each fact **510** is also associated with a source link **512** to a source of the fact.

[0067] Included in any of the facts **510** are one or more object reference links **514**. An object reference link **514** is link to object that is identified (e.g., by name) in the text of a fact. For example, in FIG. 5, in the “Wrote” fact **510c** includes an object reference link **514** to “Thriller” which is another object in the fact repository **115**. The object reference link **514** is not a hard coded link to the identified object. Instead, the object reference link **514** is a type of search link. In response to the user selecting the object reference link **514**, a search query is passed to the service engine **114**. The service engine **114** parses the search query and determines that it is a query for objects based on specified attribute and values, rather than for keywords. Accordingly, the service engine **114** passes the search request to the service engine **114**. The service engine **114** resolves the search by retrieving the objects match the search criteria.

[0068] FIG. 6 illustrates the results of the user activating the object reference link **514** in FIG. 5. Here, the results of the search by the service engine **114** returns the object “Thriller” along with its associated facts **610**, such as its writers, director, year, and running time.

Generation of Object Reference Links

[0069] As illustrated above, any fact can include one or more object reference links. In one embodiment, the object reference links are generated by one of the janitors **110**. In one embodiment, logic for generating the object reference links is follows (the implementation of course can be different, for improved efficiency).

[0070] A janitor **110** traverses the fact repository **115** and creates an objects names list comprising the name facts in the repository **115**.

[0071] Next, the janitor **110** traverses over the value portions of the facts in the repository **115** (or any portion thereof). In the value of each traversed fact, the janitor determines whether the value contains one or more of the object names on the object names list. For each object name that is found the value of the fact, the janitor creates an object reference link in the fact value, using the object name as the anchor text for the object reference link. The created object reference link is in the form of a search query, rather than a hard coded link. The query comprises one or more search criteria, which restrict the search to objects having a name that matches the object name in the name fact.

[0072] In one embodiment, the search criteria are in the form of attribute-value pairs. Accordingly, the search criteria in this embodiment is an attribute-value pair, where the attribute indicates the name fact, and the value is the name string of name contained on the object names list. An example will further illustrate this relationship.

[0073] As illustrated in FIG. 5, the “wrote” fact **510c** contained an object reference link **514** to the “Thriller (1983) (V)” object. The object reference link **514** comprises a search query of the form:

[0074] [attribute{name} value{“Thriller (1983) (V)”}].

[0075] This search query is then linked to the anchor text “Thriller (1983) (V)” in the value of fact **510c**. It would be generated in the following manner.

[0076] Upon traversing the fact repository, the janitor **110** includes the name “Thriller (1983) (V)” in the object name list. This would of course be one of hundreds, thousands, or perhaps even millions of such object names collected during a traversal. Next, the janitor, in traversing over the fact values, traverse over the facts associated with the “Michael Jackson” object, and in each of those facts **510**, processes the value portions to identify object names. In the “wrote” fact **510c**, there is the string “Thriller (1983) (V)”, and thus the janitor **110** identifies the string value “Thriller (1983) (V)”. This value is compared with the object name list, and a match is found to the “Thriller (1983) (V)” name contained therein. Accordingly, the janitor **110** creates the search query [attribute{name} value{“Thriller (1983) (V)”}] and links it to the “Thriller” string as the anchor text.

[0077] In a further embodiment, there is also the ability to trade off the precision of the search query link versus the scope of recall. This is done by including one or more additional search criteria within the query string. One implementation of this approach is including a second attribute-value pair in the search query, in which the attribute is the category fact, and the value is the value of the category fact for the object associated with the object name on the object name list. Thus, in the foregoing example, since the “Thriller” object has a category fact of “video”, the second attribute-value pair included in the query for the object reference link **514** would be

[0078] [attribute{category} value{video}].

[0079] To facilitate this refinement of the search query in the link **514**, the category attribute and value is stored along with the object name in the object name list. This approach

can be extended to include any number of additional attribute-value pairs in the search query, each one including a corresponding attribute and value from the object which is named on the object name list.

[0080] The use of search queries in the object reference links provides various advantages. Because the search query executes a general search on fact repository, it will return any new objects that are added to the fact repository after the original link was established in a previous pass by the janitor. This would not be possible if the object reference link were instead a fixed hard link to a specific, preselected fact or object.

[0081] The above process by which the janitor 110 collected object names and constructs search links in facts is preferably repeated on a periodic basis, for example daily, weekly, or any other time period or schedule. This ensures that new facts entered in the interval between repetitions are processed both to include the new object names on the object name list (so that existing facts can now properly reference such new object names), and to ensure that objects named in the value of such new facts are linked to previously existing facts. Alternatively, the process can be repeated each time some number of new facts or objects are entered (or modified) in the fact repository 115. For example, the janitor 110 may be executed after each new fact (or object), after every 100, 1000, or other threshold number of new facts (or objects) are entered (or modified).

[0082] The janitor 110 can use any type of phrase identification algorithm to identify object names in the value portions of the facts. In one embodiment, a multi-term information gain algorithm is used to identify phrases as sequences of terms in which the probability of a given sequence significantly exceeds the joint probability of the individual terms in the sequence. An implementation of this approach first tokenizes the fact value into a set of probable phrases, and then for each probable phrase checks the phrase as against the object name list to determine if the phrase is an object name. The check against the object names list is facilitated by storing the object name list as a hash table and hashing a probable phrase into the table; if the hash collides, then the phrase is an object name.

[0083] The tokenization of the terms in the value of a fact into a set of probable phrases can be implemented as follows. Each phrase will have an information gain score that is based on its probability of occurrence and the probability of occurrence of its constituent terms. A term's probability is its frequency in the corpus, per 100,000,000 terms. If the actual probability of term is unknown, it is given a default probability of 1×10^{-8} . For a sequence of terms to be considered a phrase, its probability must be higher than the joint probability of its individual terms, by a factor K , where $K=10^L$ and L is the length (number of terms) of the sequence. Thus, for a two-term sequence ($L=2$) to be considered a phrase, its actual probability must be at least $K=10^2$ times the joint probability of the individual phrases. For a three-term sequence ($L=3$) to be considered a phrase, its actual probability must be at least $K=10^3$ times the joint probability. A sequence is qualified as a phrase if its information gain satisfies this minimum ratio requirement. To facilitate rapid computation, the probabilities and scores are computed via their base 10 logarithms. Those of skill in the art can readily devise other information gain scoring algorithms.

[0084] A value of a fact then is tokenized by finding the phrases with the highest ratio of actual to expected probability. Any list traversal algorithm may be used to traverse the terms of a fact, form sequences of consecutive terms thereof, and test such sequences for being phrases using the information gain score, such as the one described above. Thus in one embodiment, sequences of terms are formed from the end (the last term) of the fact value, towards the beginning (first term), where the sequence increments in length from 1 to some number of terms. The maximum sequence length is the total number of terms between the sequence start and the first term in the value. Each such sequence is tested for its information gain score, and if the information gain is sufficient, it is deemed a phrase and memorized for later comparison against the object name list. Preferably, for each term position in the fact value, there is at most one phrase identified, being the phrase that begins at that term position that has the highest information gain. That is, if there are N terms in the fact value, there will be at most N phrases, each phrase associated with and beginning at corresponding term in the fact value. In practice, at least some of the terms of a fact value will not have any associated phrase.

[0085] The present invention has been described in particular detail with respect to one possible embodiment. Those of skill in the art will appreciate that the invention may be practiced in other embodiments. First, the particular naming of the components, capitalization of terms, the attributes, data structures, or any other programming or structural aspect is not mandatory or significant, and the mechanisms that implement the invention or its features may have different names, formats, or protocols. Further, the system may be implemented via a combination of hardware and software, as described, or entirely in hardware elements. Also, the particular division of functionality between the various system components described herein is merely exemplary, and not mandatory; functions performed by a single system component may instead be performed by multiple components, and functions performed by multiple components may instead be performed by a single component.

[0086] Some portions of above description present the features of the present invention in terms of algorithms and symbolic representations of operations on information. These algorithmic descriptions and representations are the means used by those skilled in the data processing arts to most effectively convey the substance of their invention. Furthermore, the computers referred to in the specification may include a single processor or may be architectures employing multiple processor designs for increased computing capability.

[0087] The algorithms and operations presented herein are not inherently related to any particular computer or other apparatus. Various general-purpose systems may also be used with programs in accordance with the teachings herein, or it may prove convenient to construct more specialized apparatus to perform the required method steps. The required structure for a variety of these systems will be apparent to those of skill in the art, along with equivalent variations. In addition, the present invention is not described with reference to any particular programming language. It is appreciated that a variety of programming languages may be used to implement the teachings of the present invention as

described herein, and any references to specific languages are provided for disclosure of ennoblement and best mode of the present invention.

[0088] The present invention is well suited to a wide variety of computer network systems over numerous topologies. Within this field, the configuration and management of large networks comprise storage devices and computers that are communicatively coupled to dissimilar computers and storage devices over a network, such as the Internet.

[0089] Finally, it should be noted that the language used in the specification has been principally selected for readability and instructional purposes, and may not have been selected to delineate or circumscribe the inventive subject matter. Accordingly, the disclosure of the present invention is intended to be illustrative, but not limiting, of the scope of the invention, which is set forth in the following claims. work to others skilled in the art. These operations, while described functionally or logically, are understood to be implemented by computer programs. Furthermore, it has also proven convenient at times, to refer to these arrangements of operations as modules or by functional names, without loss of generality.

[0090] Unless specifically stated otherwise as apparent from the above discussion, it is appreciated that throughout the description, discussions utilizing terms such as “processing” or “computing” or “calculating” or “determining” or “displaying” or the like, refer to the action and processes of a computer system, or similar electronic computing device, that manipulates and transforms data represented as physical (electronic) quantities within the computer system memories or registers or other such information storage, transmission or display devices.

[0091] Certain aspects of the present invention include process steps and instructions described herein in the form of an algorithm. It should be noted that the process steps and instructions of the present invention could be embodied in software, firmware or hardware, and when embodied in software, could be downloaded to reside on and be operated from different platforms used by real time network operating systems.

[0092] The present invention also relates to an apparatus for performing the operations herein. This apparatus may be specially constructed for the required purposes, or it may comprise a general-purpose computer selectively activated or reconfigured by a computer program stored on a computer readable medium that can be accessed by the computer. Such a computer program may be stored in a computer readable storage medium, such as, but is not limited to, any type of disk including floppy disks, optical disks, CD-ROMs, magnetic-optical disks, read-only memories (ROMs), random access memories (RAMs), EPROMs, EEPROMs, magnetic or optical cards, application specific integrated circuits (ASICs), or any type of media suitable for storing electronic instructions, and each coupled to a computer system bus.

We claim:

1. A computer implemented method of creating links between facts in a fact repository including a plurality of facts, each fact associated with an object, and each object including a name fact including a name of the object, the method comprising:

automatically establishing a list of object names from a plurality of name facts;

automatically identifying object names included in a plurality of facts, those object names included in the list of object names; and

for each identified object name in the plurality of facts, automatically constructing a search link in the fact, the search link comprising a search query including search criteria indicating the identified object name.

2. The method of claim 1, wherein automatically constructing a search in the fact further comprises:

using the identified object name in the fact as an anchor for the search link.

3. The method of claim 1, wherein automatically identifying, in a plurality of facts, object names, comprises:

identifying phrases in a fact; and

for each identified phrase, determining whether the phrase is included in the list of object names.

4. The method of claim 3, wherein identifying phrases in a fact comprises:

forming a plurality of sequences of consecutive terms in the fact; and

testing each sequence of terms to determine whether the sequence is a phrase.

5. The method of claim 4, wherein testing each sequence of terms comprises:

determining whether an information gain of the sequence of terms exceeds a threshold.

6. The method of claim 4, wherein testing each sequence of terms comprises:

determining whether a probability of the sequence of terms exceeds a joint probability of the terms included in the sequence.

7. The method of claim 1, further comprising:

for each object name included in the list of object name, storing a category for object associated with the object name; and

including in the search criteria of the search query an indication of the category of the object associated with the object name indicated in the search criteria.

8. The method of claim 1, further comprising:

for each object name included in the list of object names, storing an additional attribute for object associated with the object name, along with a value of the attribute; and

including in the search criteria of the search query an indication of the additional attribute and the value of the attribute of the object associated with the object name indicated in the search criteria.

9. The method of claim 1, further comprising:

displaying for an object the fact including the search link, including the object name as the anchor of the search link;

responsive to receiving a selection of the search link, providing the search query to a service engine for retrieving objects having a name fact that matches the object name included in the search criteria of the search link;

receiving from the service engine at least one search result including an object having an object name that matches the object name included in the search criteria of the search link; and

displaying the received object.

10. The method of claim 1, further comprising:

periodically repeating the steps of automatically establishing a list of object names from a plurality of name facts, automatically identifying object names included in the list of object names, and automatically constructing a search link in the fact.

11. The method of claim 1, further comprising:

responsive to a predetermined number of facts in the fact repository being modified, repeating, for each of the modified facts, the steps of automatically establishing a list of object names from a plurality of name facts, automatically identifying object names included in the list of object names, and automatically constructing a search link in the modified fact.

12. The method of claim 1, further comprising:

responsive to a predetermined number of new facts in the fact repository being created, repeating, for each of the created facts, the steps of automatically establishing a list of object names from a plurality of name facts, automatically identifying object names included in the list of object names, and automatically constructing a search link in the new fact.

13. A computer implemented method of creating links in a fact repository including a plurality of facts, each object associated with a set of facts, each fact including an attribute

descriptive of the fact and a value associated with the attribute, the values including numbers, dates, and text, the method comprising:

automatically establishing a list of fact values from a plurality of facts having a selected type of attribute;

identifying in at least one fact, a value included in the list of fact values; and

constructing a link in the fact, the link comprising a search query for an object having the identified value, and using the identified value in the fact as an anchor for the search query.

14. A user interface for a client device, for searching and browsing fact repository comprising a plurality of facts associated with objects, the fact repository coupled to a service engine adapted to search the fact repository in response to a search query, the user interface implemented on the client device by computer executable code and data, comprising:

a detailed object page, adapted for display on the client device, including a plurality of facts associated with a first object, each fact including information descriptive of the first object, at least one fact including a search link, the search link including anchor text comprising a name of a second object, and a search query for objects having a same name as the name of the second object, wherein the client device, in response to selection of the search link, is adapted to provide the search query to the service engine to retrieve the facts associated with at least one object having a same name as the name of the second object.

* * * * *