# Claim Construction Hearing for

# U.S. Patent 6,778,979

# May 19, 2011

# "selected document content"

1. A method for automatically generating a query from selected document content, comprising:

defining an organized classification of document content with each class in the organized classification of document content having associated therewith a classification label; each classification label corresponding to a category of information in an information retrieval system;

automatically identifying a set of entities in the selected document content for searching additional information related thereto using the information retrieval system;

automatically categorizing the selected document content using the organized classification of document content for assigning the selected document content a classification label from the organized classification of content; and

automatically formulating the query to restrict a search at the information retrieval system for information concerning the set of entities to the category of information in the information retrieval system identified by the assigned classification label.

# "selected document content"

## Parties' Constructions

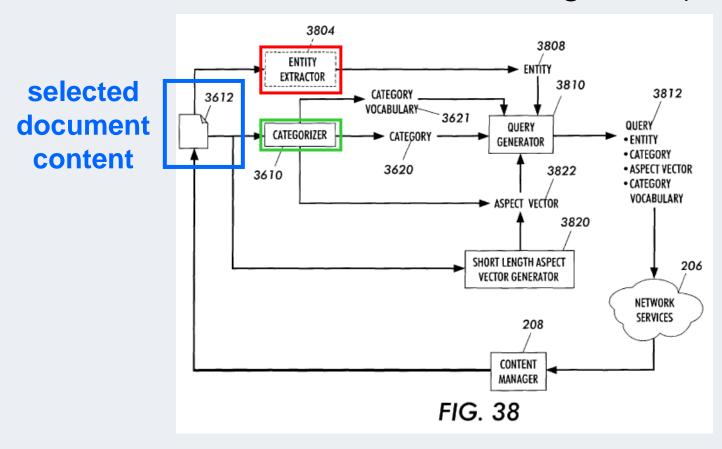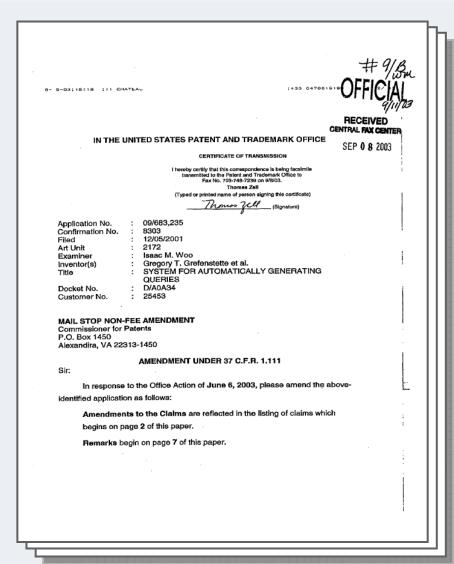| Xerox | Defendants |
|---|---|
| "all or part of the content of a document in electronic form" | Indefinite |

## Law on Indefiniteness

"A claim will be found indefinite only if it 'is insolubly ambiguous, and no narrowing construction can properly be adopted . . . .'"

Praxair, Inc. v. ATMI, Inc., 543 F.3d 1306, 1319 (Fed. Cir. 2008); Leader Techs., Inc. v. Facebook, Inc., 692 F.Supp.2d 425, 436 (D. Del. 2010).

# "selected document content"

"Selected document content" is the input to the claimed method (*i.e.*, the document content in which entities are identified and that is categorized).

**selected document content**



FIG. 38

# "selected document content"



"Applicant's claims <u>recite automatically generating a query from selected document content, from which both a set of entities and a classification label are automatically identified and assigned</u>, respectively."

9/8/2003 Applicant's Amendment at 10; emphasis in original.

# "selected document content"



**"all or part of the content of a document in electronic form"**

"In operation as shown in FIG. 38, the document content 3612 or alternatively limited context (i.e., words, sentences, or paragraphs) surrounding the entity 3808 is analyzed by categorizer 3610 to produce a set of categories 3620."

979/48:52-55; emphasis added

# "categorizing the selected document content . . ."

1. A method for automatically generating a query from selected document content, comprising:

defining an organized classification of document content with each class in the organized classification of document content having associated therewith a classification label; each classification label corresponding to a category of information in an information retrieval system;

automatically identifying a set of entities in the selected document content for searching additional information related thereto using the information retrieval system;

automatically categorizing the selected document content using the organized classification of document content for assigning the selected document content a classification label from the organized classification of content; and

automatically formulating the query to restrict a search at the information retrieval system for information concerning the set of entities to the category of information in the information retrieval system identified by the assigned classification label.

# "categorizing the selected document content . . ."

## Parties' Constructions

| Xerox | Defendants |
|---|---|
| "determining the subject matter of the selected document content using one or more of the categories defining the organized classification of document content and assigning the corresponding classification label(s) to the selected document content." | "using the organized classification of document content to categorize the selected document content and to assign to the selected document content a single classification label." |

# "categorizing the selected document content . . ."

## The "crux of the . . . dispute" according to Defendants

"The crux of the parties' dispute is whether, as Defendants contend, a single classification label is assigned in the 'categorizing' step and used to identify 'the' single category used to restrict a search in the 'formulating' step, or whether, as Xerox contends, more than one classification label may be assigned in the 'categorizing' step and employed in the 'formulating' step."

Defendants' Opening Br. at 11.

# "categorizing the selected document content . . ."

## Law on "a" meaning "one or more"

"That 'a' or 'an' can mean 'one or more' is best described as a rule, rather than merely as a presumption or even a convention. The exceptions to this rule are extremely limited: a patentee must evince a clear intent to limit 'a' or 'an' to 'one'."

Baldwin Graphic Sys., Inc. v. Siebert, Inc., 512 F.3d 1338, 1342-43 (Fed. Cir. 2008) (internal quotations omitted).

# "categorizing the selected document content . . ."

## Use of one or more labels + categories

1. A method for automatically generating a query from selected document content, comprising:

defining an organized classification of document content with each class in the organized classification of document content having associated therewith a classification label; each classification label corresponding to a category of information in an information retrieval system;

automatically identifying a set of entities in the selected document content for searching additional information related thereto using the information retrieval system;

automatically categorizing the selected document content using the organized classification of document content for assigning the selected document content a classification label from the organized classification of content; and

automatically formulating the query to restrict a search at the information retrieval system for information concerning the set of entities to the category of information in the information retrieval system identified by the assigned classification label.

(1) "a" means "one or more"

(2) one or more labels correspond to one or more categories of information in IRS

(3) "the" refers back to one or more labels / categories

(4) search is restricted to one or more categories of information in IRS (identified by one or more labels)

# "categorizing the selected document content . . ."

**The '979 Patent teaches the use of one or more labels.**

"Document <u>classification labels</u> define the <u>set of categories</u> 3620 output by the categorizer 3610. These <u>classification labels</u> in one embodiment are appended to the query 3812 by query generator 3810 to restrict the scope of the query (i.e., the entity 3808 and the context vector 3822) to folders corresponding to <u>classification labels</u> in a document collection of an information retrieval system."

979/49:31-37; emphasis added

# "classification label"

1. A method for automatically generating a query from selected document content, comprising:

defining an organized classification of document content with each class in the organized classification of document content having associated therewith a classification label; each classification label corresponding to a category of information in an information retrieval system;

automatically identifying a set of entities in the selected document content for searching additional information related thereto using the information retrieval system;

automatically categorizing the selected document content using the organized classification of document content for assigning the selected document content a classification label from the organized classification of content; and

automatically formulating the query to restrict a search at the information retrieval system for information concerning the set of entities to the category of information in the information retrieval system identified by the assigned classification label.

# "classification label"

## Parties' Constructions

| Xerox | Defendants |
|---|---|
| "a label in any format that identifies a category in the organized classification of document content." | "classifying word or phrase" |

# "classification label"

## Example of Defendants' General Dictionary Definitions

**¹la·bel** \ˈlābəl\ *n* -s *often attrib* [ME fr. MF, fr. OF *label* ribbon, fringe, label in heraldry, prob. of Gmc origin; akin to OHG *lappa* flap, lappet — more at LAP] **1** *archaic* : a narrow

. . .

book) ⟨read the ~ on the bottle⟩ **c** : a descriptive, classifying, or identifying word or phrase: as (1) : EPITHET ⟨the term stream of consciousness . . . is already established as a literary ~ —Robert Humphrey⟩ ⟨acquired the ~ of "playboy" which seemed to stick —Brian Crozier⟩ ⟨hanging the subversive ~ on their own liberal clergy —Ralph Winnett⟩ (2) : a word or phrase used with but not as part of a dictionary definition usu. in abbreviated form and distinctive type to provide information (as grammatical function or area or level of usage) about the word defined ⟨the ~ *obsolete* is abbreviated *obs*⟩ (3) : a newspaper headline merely identifying the subject matter of an article rather than summarizing action **6** : a

Webster's Third New Int'l Dictionary, Unabridged (2002)

# "classification label"

## Law on Use of General Dictionaries

"By design, general dictionaries collect the definitions of a term as used not only in a particular art field, but in many different settings. . . . For that reason, we have stated that a general-usage dictionary cannot overcome art-specific evidence of the meaning of a claim term."

Phillips v. AWH Corp., 415 F.3d 1303, 1321-22 (Fed. Cir. 2005).

# "classification label"

## Examples of Art-Specific Dictionaries

"[a]n identifier within or attached to a set of data elements."

IBM Dictionary of Computing (1994).

"[a] data item that serves to identify a data record (much in the same way as a key is used) . . . ."

McGraw-Hill Dictionary of Computing & Communications (2003)

# "query"

1. A method for automatically generating a <mark>query</mark> from selected document content, comprising:

defining an organized classification of document content with each class in the organized classification of document content having associated therewith a classification label; each classification label corresponding to a category of information in an information retrieval system;

automatically identifying a set of entities in the selected document content for searching additional information related thereto using the information retrieval system;

automatically categorizing the selected document content using the organized classification of document content for assigning the selected document content a classification label from the organized classification of content; and

automatically formulating the <mark>query</mark> to restrict a search at the information retrieval system for information concerning the set of entities to the category of information in the information retrieval system identified by the assigned classification label.

# "query"

## Parties' Constructions

| Xerox | Defendants |
|---|---|
| "a set of data specifying search criteria" | "request for search results" |

# "query"

**1.** A method for automatically generating a query from selected document content, comprising:

defining an organized classification of document content with each class in the organized classification of document content having associated therewith a classification label; each classification label corresponding to a category of information in an information retrieval system;

automatically identifying a set of entities in the selected document content for searching additional information related thereto using the information retrieval system;

automatically categorizing the selected document content using the organized classification of document content for assigning the selected document content a classification label from the organized classification of content; and

automatically formulating the query to restrict a search at the information retrieval system for information concerning the set of entities to the category of information in the information retrieval system identified by the assigned classification label.

# "query"

A "query" in the '979 Patent is defined by its contents



"The query generated may include some or all of the following elements…: (a) a set of entities 3808…, (b) a set of categories 3620 generated by the categorizer 3610…, (c) an aspect vector 3822 generated by categorizer 3610 or short run aspect vector generator 3820, and (d) a category vocabulary 3621 generated by the categorizer 3610"

979/48:41-51

# "to restrict . . . label"

1. A method for automatically generating a query from selected document content, comprising:

defining an organized classification of document content with each class in the organized classification of document content having associated therewith a classification label; each classification label corresponding to a category of information in an information retrieval system;

automatically identifying a set of entities in the selected document content for searching additional information related thereto using the information retrieval system;

automatically categorizing the selected document content using the organized classification of document content for assigning the selected document content a classification label from the organized classification of content; and

automatically formulating the query to restrict a search at the information retrieval system for information concerning the set of entities to the category of information in the information retrieval system identified by the assigned classification label.

# "to restrict . . . label"

## Parties' Constructions

| Xerox | Defendants |
|---|---|
| "the set of data specifying search criteria includes data items corresponding to one or more entities identified in the 'automatically identifying' step and one or more classification labels assigned in the 'automatically categorizing' step." | "to confine a search at the information retrieval system to the category of information identified by the assigned classification label, where the search seeks information concerning the set of entities." |

# "to restrict . . . label"

1. A method for automatically generating a query from selected document content, comprising:

defining an organized classification of document content with each class in the organized classification of document content having associated therewith a classification label; each classification label corresponding to a category of information in an information retrieval system;

automatically identifying a set of entities in the selected document content for searching additional information related thereto using the information retrieval system;

automatically categorizing the selected document content using the organized classification of document content for assigning the selected document content a classification label from the organized classification of content; and

automatically formulating the query to restrict a search at the information retrieval system for information concerning the set of entities to the category of information in the information retrieval system identified by the assigned classification label.

# "to restrict . . . label"

The specification portion Defendants cite confirms that Xerox's construction is correct, despite Defendants' misleading underlining.

> "The specification explains that 'the search is focused on documents found in the single node of the document hierarchy genetics, at 3910'."

Defs.' Br. at 8 (quoting 979/50:10-11);
emphasis in original.

# Order of Steps (Claims 1, 18)

## Parties' Constructions (Claims 1, 18)

| <u>Xerox</u> | <u>Defendants</u> |
|---|---|
| Step (a) must be performed before steps (c) and (d). | Step (a) must be performed before steps (c) and (d). |
| Step (b) must be performed before the completion of step (d). | Step (b) must be performed before step (d). |
| Step (c) must be performed before the  completion of step (d). | Step (c) must be performed before step (d). |

# Order of Steps (Claims 1, 18)

## Law on order of steps

"[A]s a general rule the claim is not limited to performance of the steps in the order recited, unless the claim explicitly or implicitly _requires_ a specific order."

Baldwin Graphic Sys., Inc. v. Siebert, Inc., 512 F.3d 1338, 1345 (Fed. Cir. 2008) (emphasis added).

# Order of Steps (Claims 1, 18)

## Law on order of steps

"First, we look to the claim language to determine if, as a matter of logic or grammar, they must be performed in the order written. . . . If not, we next look to the rest of the specification to determine whether it directly or implicitly requires such a narrow construction. If not, the sequence in which such steps are written is not a requirement."

Altiris, Inc. v. Symantec Corp., 318 F.3d 1363, 1369-70 (Fed. Cir. 2003) (internal citations omitted).

# Order of Steps (Claims 1, 18)

"The computer system may be implemented by any one of a plurality of configurations. For example, processor may in alternative embodiments, be defined by a collection of microprocessors configured for multiprocessing. In yet other embodiments, the functions provided by software components may be distributed across multiple computing devices (such as computers and peripheral devices) acting together as a single processing unit."

979/75:15-22

## Parties' Constructions (Claims 1 & 2)

| Xerox | Defendants |
|---|---|
| The step of Claim 2 must be performed during or after the completion of step (d) of Claim 1. | The steps of claim 1 must be performed before the step of 2. |

# Order of Steps (Claims 1 & 2; 18 & 19)

**"Entities" (Claim 1) and "terms" (Claim 2) can be present in the formulated query <u>before</u> "categories" (Claim 1).**

"Document classification labels define the set of categories 3620 output by the categorizer 3610. These classification labels in one embodiment are <u>appended</u> to the query 3812 by query generator 3810 to restrict the scope of the query (i.e., the <u>entity</u> 3808 and the <u>context vector</u> 3822)"

979/49:31-35; emphasis added.

# The Parties' Deposition History

Apr. 25/May 10: Xerox requests follow-up to Google 30(b)(6)

Yesterday, Google offered June 14 for this deposition

**?**

**5-7 wks**

Apr. 23: Defs request Chuat/Fernstrom depositions

June 29/July 1: Xerox's proposed deposition dates

**9 wks +**

Feb. 7: Xerox requests 30(b)(6) depositions of Yahoo

May 3-6: Yahoo 30(b)(6) depositions occur

**12 wks +**

Feb. 1: Xerox requests 30(b)(6) deposition of Google

Apr. 7: Google 30(b)(6) deposition occurs

**9 wks +**

Jan. 28: Defs request depositions of inventors / Zell

Mar. 30, Apr. 1 & 7: Inventor / Zell depositions occur

**9 wks +**

| Jan | Feb | Mar | Apr | May | June |