

(19)  
(12)

(KR)  
(A)

(51) 。 Int. Cl. <sup>7</sup>  
G06F 17/30

(11)  
(43)

2002 - 0003701  
2002 01 15

(21) 10 - 2000 - 0035714  
(22) 2000 06 27

(71)  
4 563 6  
- - 133 2 7

(72)  
- - 133 27

(74)  
:

(54)

(skip)

3

1 ;

- 2 ;
- 3 ;
- 4 ;
- 5 ;
- 6A 6B ;
- 7A 7C 가 ;
- 8 OCR ;
- 9 OCR ;
- 10A 10B ;
- 11 .

(indexing)

(Internet)

(retrieval system)

mparison),

(automatic classification),

(automatic indexing),  
(relevant feedback)

(approximate co

가

(OCR(Optical Character Recognition)

MIDI(Musical Instrument Digital Interface)

가  
가  
가  
가

가  
가  
가  
가

가  
가  
가  
가

가

" (vocabulary mismatch)" 가

(key melody)

( ) (phrase database) 가

가 가

, OCR

1 가 (pseudo code) 가

가 ( )

(book index) 가

39%), 63% , (19%) (42%), (

, PAT - 81% . PAT - 가

600 가 (Magabytes)  
가

PAT -

OCR (noisy documents)

(layout principles)

(newborn phrases)

가

가

ion) 가 가 (repetit  
가 가

(string)

가

가 ) (Hsu, Liu, Chen, 1998).

$O(n^2, n)$   
(Hsu, Liu, Chen, 1999)

가  $O(n \log n)$  , ABCABCAABCAABCABCABCA  
ABCABCA가

1972 Karp, Miller & Rosenberg  
가 (1) 가 K  
가

1995 Soldano, Viari & Champesme  
(2) 가 L  
가 가

BCA 3  
BCA ABCA  
ABCA가 가 4  
EFG

ABC, BCA  
ABCA

EFGABCABCAEFG  
ABC BCA ABCA  
가 가 가  
EFG, ABC, ABC,  
ABCA EFG  
ABCA

1. 가 , 가  
(3) (1) (2)

2. 가,  
OCR

3. ,

4.  $O(n^2)$  ,

5. ( 가 L ) 가 K

가 (separator) 가 (reference mark) 가

가 (skip) (taken out) 가  
가 가  
가 가  
가 가  
가 가

1. , , , .
  2. , , OCR , (noisy document)
  3. , , .
  4. , , 가 , .
  5. (deterministic principles) , ( , 가 2 ) 가 .
  6. , , , , , DNA , 가 .
  7. (computing complexity)  $O(L*N)$  , L 가 , N .
  8. ( , DVD , , , , , ) ( )
- 2 11 .
- , " E F G A B C A B C A E F G " , , , , DNA .
- 가 (element) , (combine), (erase) 가 (accept) , 가 .
- , (stop condition) (sort) (filtering) . 가 ( , , ) .

2 가 (pseudo code) , 가 Freq()

가 (hash function)

, O(1) 가 20% 가 14%

(position list) 가

3 (:) (listed elements) 3 가

(listed data structure) (procedure) , x (separator)

가

가 (taken - out) 가

(combination condition) (acceptance condition) (erasure condition)

( , Freq(LIST[I]) Freq(LIST[I+1])) 1

가 1

가 가

(global computation) 가 (local computation) 1 2

가

1 가 , O  
 (N)가 2 , N ,  
 가 , 가 , 가 O(N)  
 가 L L  
 가 L , 2 O(L\*N) O(1)가 가

가 가

4 (position list) 4  
 (combination)  
 (E, F, G, A, B, C, A, B, C, A, E, F, G) (1 11, 2, 12, 3 13, 4 7 10, 5 8, 6 9, 4 7  
 10, 5 8, 6 9, 1 11, 2, 12, 3, 13) E 1 11 F 2 12  
 3 가 ( , EFG:2, x, ABC:2, BCA:2, CAB:1, ABC:2, BCA:2, x, EFG  
 :2) (1 11, x, 4 7, 5 8, 6, 4 7, 5 8, x, 1 11)

1  
(local decrement process)

가 가 ,  
 4 " LIST" " 4 7 10" " 5 8" " A" " B"  
 " A" 4, 7 10 " B"  
 5 8 " A" " B" 4 5 7 8  
 " A" " B" 4 7 " 4 7"  
 , " A" " B" 가

가  
(loop)가

가 1 가 , 가  
 " 1 11" " EFG" , 가 3 가

가 1 가 , 가  
 가 가  
 가

2 ~ 3 , 가 20%

, 가 14% . 가 가 2  
2 3 40% . 2N 3N  
( , O(N)) .

5 가 가 가 가 가  
가 가 O(LN) 가 가 가

5 , 2 ,  
3 .

(stop words;  
, " of" , " on" , " the" ) .  
가 " 的" 가 6A 6B  
(  
) . " FGDC 的 Digital Geospatial metadata" .

(fetch) , 2- 가 127  
가 가 가 (punctuation mark) 가 가  
(legal characters) (1)  
(2)  
가

7A 7C . 7C 7B  
가 가 ,

. 35 11 , 가 8982  
859 , 66500 phase 90.4% 13035 8123 ,  
57389

, 911 . 86.3% .

(term suggestion)

(term relevance feedback)

8 OCR

40709 , 8 , , 1300 OCR  
 , " 鄉村幹部" , " 農村幹部" , " 鄉村幹都" , " 工農幹部" " 鄉村幹部( )" 가  
 , 가 (match) . OCR  
 OCR , 가 . ,

9 OCR  
 OCR , 가

10B " DVD" 가 , " DVD" 가 482 10A  
 가 , " DVD player"

11 " " " 17 - 1321" 가  
 (Ave Maria) , " " .  
 , 가 , . ,

(relevance feedback),

가

(57)

1.

:

ference mark)      가      ,      (separator)      (re

(A)      ;

(B)      ;

(C)      ;

가      ,      ;

가      가      ,      가  
가      ;

가      가      ,      가  
가      가      ;

가      ,      가      ;

(D)      ;      (B) - (C)

;

(E)

2.

1      ,

가      ,      가

3.

1      ,

가 2

4.

1 ,

5.

4 ,

6.

1 ,

7.

1 ,

8.

7 ,

9.

1 ,

가

10.

9 ,

가  
가

가

가

11.

10

가

가

12.

11

13.

12

14.

1

15.

1

가

16.

1

17.

16

DNA



1. Convert the input into a *LIST*.
2. Do Loop
  - 2.1 Set *MergeList* to *empty*.
  - 2.2 Put a *separator* to the end of *LIST* as a sentinel and set the occurring frequency of the *separator* to 0.
  - 2.3 For *I* from 1 to  $\text{NumOf}(\text{LIST}) - 1$  step 1, do
    - If  $\text{LIST}[I]$  is the *separator*, Go to Label 2.3.
    - If  $\text{Freq}(\text{LIST}[I]) > \text{threshold}$  and  $\text{Freq}(\text{LIST}[I+1]) > \text{threshold}$ 
      - then
        - Merge  $\text{LIST}[I]$  and  $\text{LIST}[I+1]$  into *Z*
        - Put *Z* to the end of *MergeList*.
      - else
        - If  $\text{Freq}(\text{LIST}[I]) > \text{threshold}$  and  $\text{LIST}[I]$  did not merge with  $\text{LIST}[I-1]$ 
          - then
            - Save  $\text{LIST}[I]$  in *FinalList*.
            - If the last element of *MergeList* is not the *separator*,
              - Then
                - Put the *separator* to the end of *MergeList*.
- 2.4 Set *LIST* to *MergeList*.
- Until  $\text{NumOf}(\text{LIST}) < 2$ .

Example 1: Given : E F G A B C A B C A E F G. Let threshold=1,  
separator=x

Step 1: Create a list of single characters:

*LIST* = (E:2, F:2, G:2, A:3, B:2, C:2, A:3, B:2, C:2, A:3, E:2, F:2, G:2)

Step 2:

After 1st iteration :

*MergeList* = (EF:2, FG:2, GA:1, AB:2, BC:2, CA:2, AB:2, BC:2,  
CA:2, AE:1, EF:2, FG:2)

*FinalList* = ( )

After 2nd iteration :

*MergeList* = (EFG:2, x, ABC:2, BCA:2, CAB:1, ABC:2, BCA:2, x,  
EFG:2)

*FinalList* = ( )

After 3rd iteration :

*MergeList* = (x, ABCA:2, x, ABCA:2)

*FinalList* = (EFG:2 )

After 4th iteration :

*MergeList* = (x)

*FinalList* = (EFG:2, ABCA:2)

Example 2: Given : E F G A B C A B C A E F G. Let threshold=1,  
separator=x

Step 1: Create a position list :

$LIST = (1\ 11, 2\ 12, 3\ 13, 4\ 7\ 10, 5\ 8, 6\ 9, 4\ 7\ 10, 5\ 8, 6\ 9, 4\ 7\ 10, 1\ 11, 2\ 12, 3\ 13)$

Step 2:

After 1st iteration :

$MergeList = (1\ 11, 2\ 12, 3, 4\ 7, 5\ 8, 6\ 9, 4\ 7, 5\ 8, 6\ 9, 10, 1\ 11, 2\ 12)$

$FinalList = ()$

After 2nd iteration :

$MergeList = (1\ 11, x, 4\ 7, 5\ 8, 6, 4\ 7, 5\ 8, x, 1\ 11)$

$FinalList = ()$

After 3rd iteration :

$MergeList = (x, 4\ 7, x, 4\ 7)$

$FinalList = (1\ 11) = (E\ F\ G : 2)$

After 4th iteration :

$MergeList = (x)$

$FinalList = (1\ 11, 4\ 7) = (E\ F\ G : 2, A\ B\ C\ A : 2)$

1. Convert the input string into a *LIST* of single characters and accumulate their occurring frequencies.
2. Do Loop
  - 2.1 Set *MergeList* to *empty*.
  - 2.2 Put a *separator* to the end of *LIST* as a sentinel and set the occurring frequency of the *separator* to 0.  
If the length of the element in *LIST* is less than *SomeNumber*,  
Then  
    Implement step 2.3 by use of hash functions.  
Else  
    Implement step 2.3 with the position list method.
  - 2.4 Set *LIST* to *MergeList*.  
Until  $\text{NumOf}(\text{LIST}) < 2$ .
3. Filter the *FinalList* and sort the result according to some criteria.

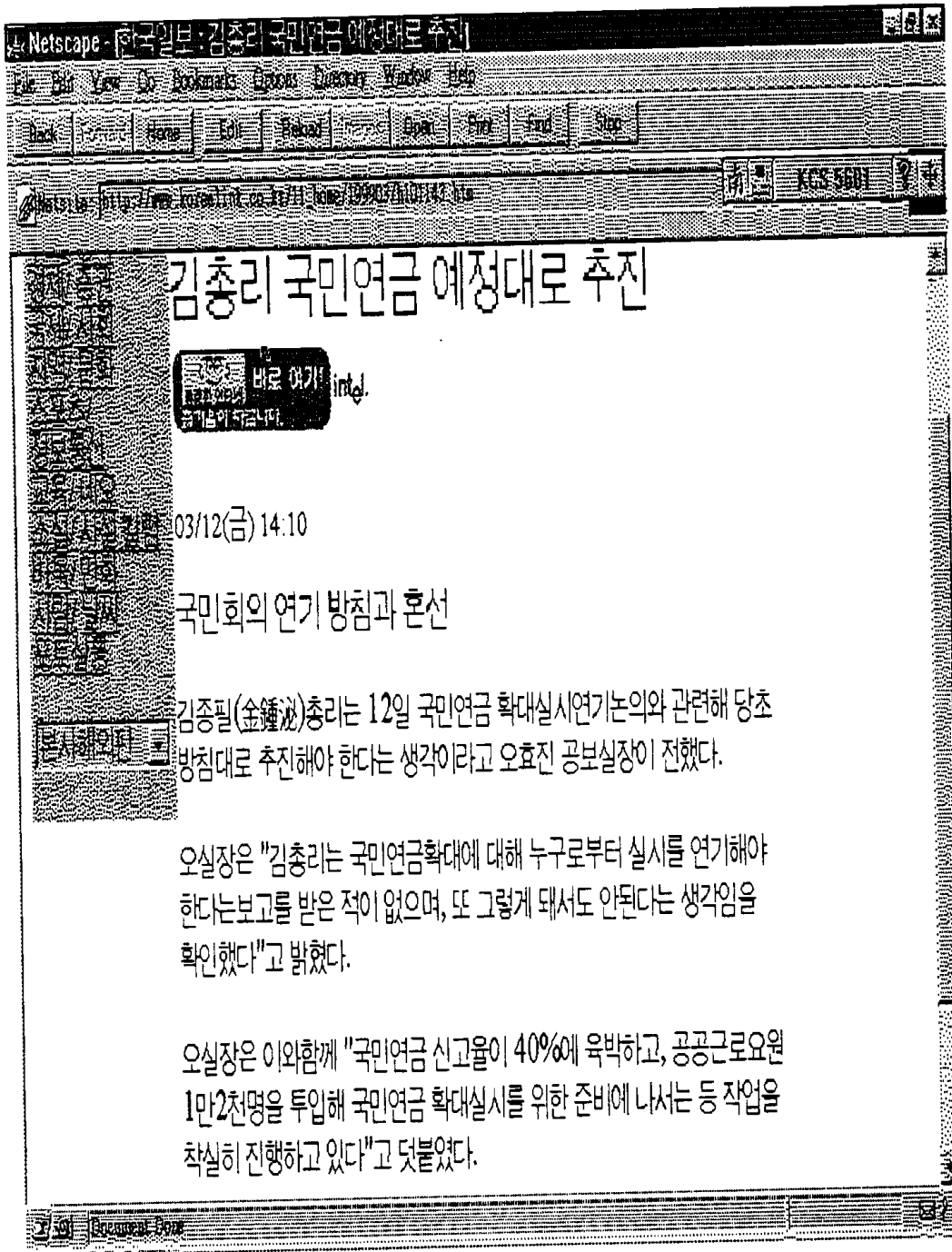
### **Comparison of Three Metadata Related Standards**

在本文中，我們介紹了三個跟 metadata 相關的標準，它們分別是 FGDC 的 Digital Geospatial Metadata、Dublin Core、和 URC。雖然它們各有自己的設計目標和特質，但都是假設其操作環境為類似網際網路的環境。FGDC 的 Digital Geospatial Metadata 是設計來專門處理地理性資料，由於它有聯邦行政命令的支持，可說是已成為美國在地理方面的資料著錄國家標準。Dublin Core 則比較像是 USMARC 的網路節縮版，使非專業人士也能在短時間內熟悉和使用此格式來著錄收藏資料，但在現階段祇針對類似傳統印刷品的電子文件。由 IETE 的 URI 工作小組所負責的 URC，其原始的設計目的雖是用來連結 URL 和 URN，但為因應電子圖書館時代的要求，其內含逐漸擴大，雖然尚在發展中，但由於有 IETE 的支持，未來成為網際網路上通用標準的可能性極大。在此文中，我們也從幾個不同角度，分析和比較這三個 metadata 格式的異同和優缺點。

6b

Terms Before Filtering		Terms After Filtering	
1.	設計 : 3 (design)	1.	設計 : 3 (design)
2.	資料 : 3 (data)	2.	資料 : 3 (data)
3.	網路 : 3 (network)	3.	網路 : 3 (network)
4.	標準 : 3 (standard)	4.	標準 : 3 (standard)
5.	Dublin Core : 2	5.	Dublin Core : 2
6.	FGDC 的 Digital Geospatial Metadata : 2	6.	FGDC 的 Digital Geospatial Metadata : 2
7.	IETE的 : 2 (of IETE)	7.	IETE : 2
8.	三個 : 2 (three)	8.	三個 : 2 (three)
9.	文中 : 2 (in the article)	9.	文中 : 2 (in the article)
10.	比較 : 2 (comparison)	10.	比較 : 2 (comparison)
11.	它們 : 2 (they)	11.	它們 : 2 (they)
12.	由於 : 2 (owing to)	12.	由於 : 2 (owing to)
13.	地理 : 2 (geography)	13.	地理 : 2 (geography)
14.	成爲 : 2 (become)	14.	成爲 : 2 (become)
15.	我們 : 2 (we)	15.	我們 : 2 (we)
16.	的支持 : 2 (support of)	16.	支持 : 2 (support)
17.	的設計目 : 2 (incorrect term)	17.	設計目 : 2 (incorrect term)
18.	格式 : 2 (format)	18.	格式 : 2 (format)
19.	記錄 : 2 (record)	19.	記錄 : 2 (record)
20.	電子 : 2 (electronics)	20.	電子 : 2 (electronics)
21.	網際網路 : 2 (Internet)	21.	網際網路 : 2 (Internet)
22.	環境 : 2 (environment)	22.	環境 : 2 (environment)
23.	雖然 : 2 (although)	23.	雖然 : 2 (although)
24.	類似 : 2 (similar)	24.	類似 : 2 (similar)

7a



7b

Netscape - [Results of Automatic Keyword Extraction]

File Edit View Go Bookmarks Custom Directory Window Help

Back Forward Home Stop Reload Open Print Find

http://www.kci.go.kr/web/keyword/extract.asp

Sort by term frequency	Sort by term length	Broader and Narrower Terms
1. 국민:6	1. 국민연금확대실시:2	1. 국민:6
2. 국민연금:5	2. 국민연금확대:3	1. 국민연금:5
3. 국민연금확대:3	3. 해야한다는:2	2. 국민연금확대:3
4. 실장:3	4. 국민연금:5	3. 국민연금확대실시:2
5. 연기:3	5. 다는생각:2	2. 실장:3
6. 국민연금확대실시:2	6. 대로추진:2	1. 오실장은:2
7. 김총리:2	7. 오실장은:2	3. 연기:3
8. 다는생각:2	8. 김총리:2	4. 김총리:2
9. 대로추진:2	9. 실시를:2	5. 다는생각:2
10. 방침:2	10. 총리는:2	6. 대로추진:2
11. 실시를:2	11. 국민:6	7. 방침:2
12. 오실장은:2	12. 실장:3	8. 실시를:2
13. 총리는:2	13. 연기:3	9. 총리는:2
14. 하고:2	14. 방침:2	10. 하고:2
15. 해야한다는:2	15. 하고:2	11. 해야한다는:2
16. 했다:2	16. 했다:2	12. 했다:2

7c

## President Kin says: National Annuity is Encouraged to be Executed

1. National Annuity is Enlarged to be Executed	9. Execution
2. National Annuity is Enlarged	10. President
3. Must be Done	11. National
4. National Annuity	12. Action Officer
5. Thoughts of	13. Delayed
6. Increase	14. Guiding Policy
7. Action Officer Wu	15. Do
8. President Kin	16. Complete

Search Results for 鄉村幹部 - Netscape

File Edit View Go Communicator Help

Among 40709 records, find 120 records similar to  
「鄉村幹部」

match score	keywords	term frequency
1000	<input checked="" type="checkbox"/> 鄉村幹部	4
695	<input type="checkbox"/> 村幹部	5
695	<input checked="" type="checkbox"/> 農村幹部	4
695	<input type="checkbox"/> 一些鄉村幹	2
695	<input checked="" type="checkbox"/> 鄉村幹部	2
391	<input type="checkbox"/> 幹部	827
391	<input type="checkbox"/> 鄉村	46
391	<input checked="" type="checkbox"/> 農幹部	7
391	<input type="checkbox"/> 優秀幹部	7
391	<input type="checkbox"/> 名幹部	7

0 20 40 60 80 100

Or add other terms

Select:  Keyword database  Document database

開始查詢 / Query

Search Results for 反革命份子 - Netscape

File Edit View Go Communicator Help

Among 40709 records, find 130 records similar to 「反革命份子」

match score	keywords	term frequency
1000	<input checked="" type="checkbox"/> <u>反革命份子</u>	3
600	<input checked="" type="checkbox"/> <u>反革命分子</u>	258
600	<input type="checkbox"/> <u>個反革命分子</u>	8
600	<input type="checkbox"/> <u>一切反革命分子</u>	6
600	<input type="checkbox"/> <u>清一切反革命分子</u>	5
600	<input checked="" type="checkbox"/> <u>反革命分子萬玉昌</u>	4
600	<input type="checkbox"/> <u>名反革命分子</u>	3
600	<input type="checkbox"/> <u>切反革命分子</u>	3
600	<input type="checkbox"/> <u>反革命政治煽子</u>	3
600	<input type="checkbox"/> <u>批反革命分子</u>	3
600	<input type="checkbox"/> <u>首的反革命分子</u>	2
600	<input type="checkbox"/> <u>反革命分子王</u>	2
600	<input checked="" type="checkbox"/> <u>反革命分子王洪富</u>	2
600	<input type="checkbox"/> <u>反革命分子的材料</u>	2
600	<input type="checkbox"/> <u>包庇反革命分子</u>	2

10a

Search Results for DVD - Netscape

File Edit View Go Communicator Help

Among 58242 records, find 297 records similar to 「DVD」

match score	keywords	term frequency
1000	<input type="checkbox"/> <u>DVD</u>	482
1000	<input checked="" type="checkbox"/> <u>DVD播放機</u>	89
1000	<input checked="" type="checkbox"/> <u>DVD Player</u>	76
1000	<input type="checkbox"/> <u>DVD Audio</u>	19
1000	<input type="checkbox"/> <u>DVD光碟機</u>	10
1000	<input type="checkbox"/> <u>DVD播放機市場</u>	9
1000	<input type="checkbox"/> <u>DVD Forum</u>	8
1000	<input type="checkbox"/> <u>DVD軟體</u>	8
1000	<input type="checkbox"/> <u>DVD市場</u>	7
1000	<input type="checkbox"/> <u>DVD產品</u>	6

10b

Search Results for DVD Player, DVD播放...

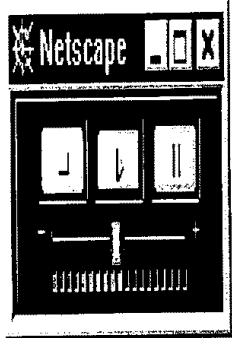
File Edit View Go Communicator Help

Among 58242 records, find 297 records similar to  
「DVD Player, DVD播放機」

match score	keywords	term frequency
781	<input type="checkbox"/> <u>DVD Player機種</u>	2
781	<input type="checkbox"/> <u>DVD Player機體模組</u>	1
781	<input checked="" type="checkbox"/> <u>DVD Player二代機</u>	1
750	<input type="checkbox"/> <u>DVD播放機</u>	89
750	<input type="checkbox"/> <u>DVD Player</u>	76
750	<input type="checkbox"/> <u>DVD播放機市場</u>	9
750	<input type="checkbox"/> <u>生產 DVD播放機</u>	6
750	<input checked="" type="checkbox"/> <u>可攜式 DVD播放機</u>	6
750	<input type="checkbox"/> <u>產 DVD Player</u>	4
750	<input type="checkbox"/> <u>DVD Player的生產據點</u>	3

Among 444 records, find 62 records similar to  
 play

Query Candidates:



Match Score	Title / Track / Theme No.	Theme (length)	Theme Frequency
	<a href="#">click to find similar titles or tracks</a>	<a href="#">click to find similar music</a>	<a href="#">click to play</a>
1: 1000	<input type="checkbox"/> <a href="#">Schubert: Ave Maria // 0</a>	<input type="checkbox"/> <a href="#">As5 A5 As5 D6 C6 As5</a> (5)	2
2: 838	<input type="checkbox"/> <a href="#">BEETHOVEN: Symphony No. 4, 3rd. Mvt. / Viola / 0</a>	<input type="checkbox"/> <a href="#">Cs5 Cs5 Cs5 Cs5 Cs5 ...</a> (711)	2
3: 716	<input type="checkbox"/> <a href="#">BEETHOVEN: Piano Sonata No. 4 // 0</a>	<input type="checkbox"/> <a href="#">As5 G5 Ds5 G5 As5 Ds6 ...</a> (252)	2

Others : [0](#) [3](#) [6](#) [9](#) [12](#) [15](#) [18](#) [21](#) [24](#) [27](#) (30 After)

Or add other terms

Select:  Key melodies  Music database