



Chinese World

Tsuyi

編碼標準

□ ASCII

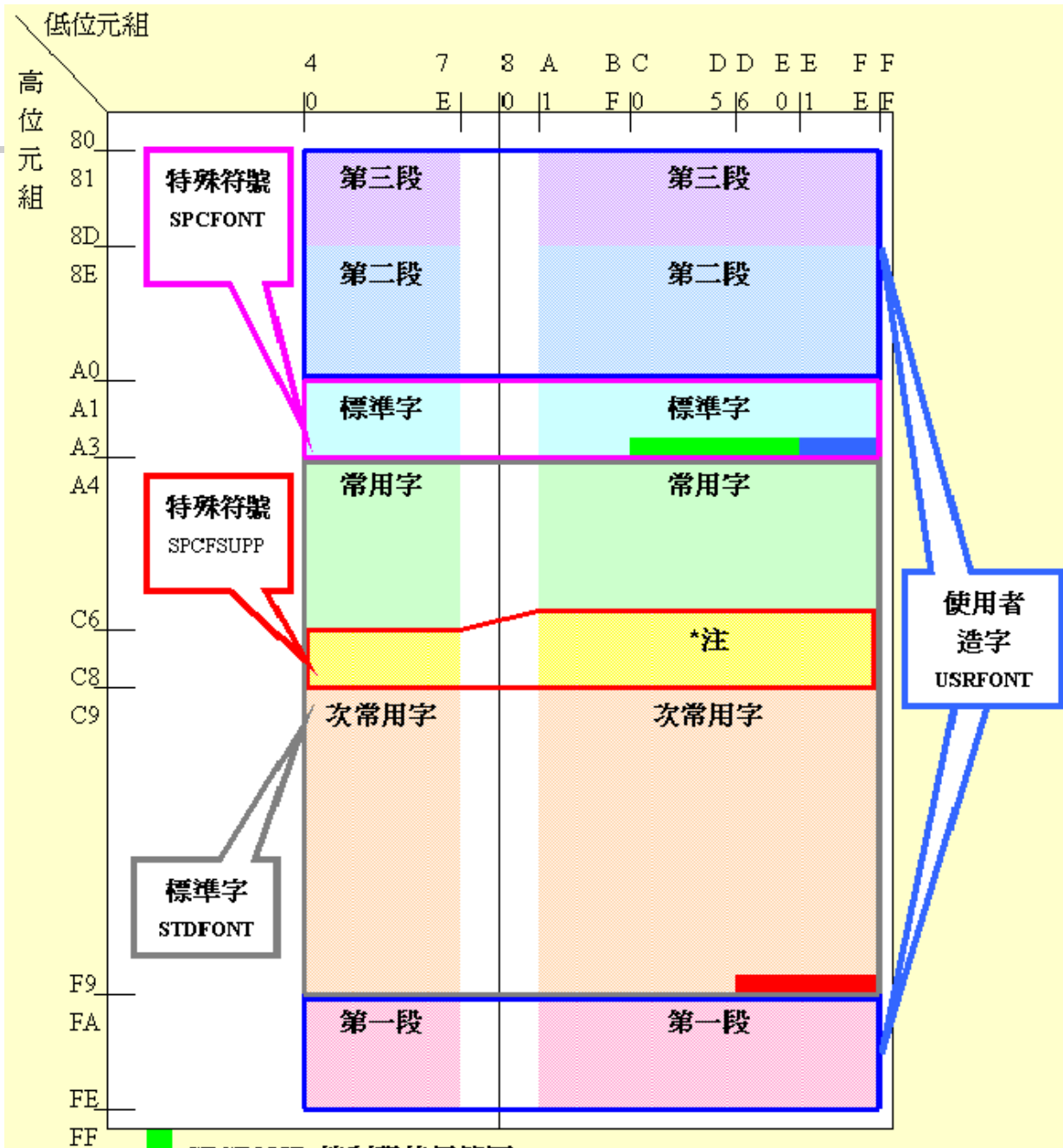
- 8 bits (理論上有 256 種可能)
- 0x00 ~ 0x7F 共 128 種字元
 - 0x00 ~ 0x1F → control characters
 - 0x20 ~ 0x7F → printable characters

□ Big5

- 使用 2 bytes 來存放中文字 (理論上有 65536 種可能)
- 實際上為與 ASCII 相容, 只能使用 19782 個
 - $[0x81 \sim 0xFE][0x40 \sim 0x7E, 0xA1 \sim 0xFE]$
 $= 126 * (63 + 94) = 126 * 157 = 19782$

編碼標準 - Big5

- 標準字 (13502)
 - 常用字
 - 你我他的媽
 - 次常用字
 - 杓晃束鏢廳
- 特殊符號 (441)
 - 符號、控制碼
 - : ! ° ∩ † ‡
 - 罕用符號
- 使用者造字 (5809)
 - 三段



Big5的問題

□ 使用者造字區

- 每個人都可以自己造字
於是自己造的字放到別人電腦上就看不到

□ 延伸版本繁雜

- 倚天Big5延伸
- Code Page 950
- Big5+
- 族繁不及備載..

□ 許功蓋 (\)

- 許 (0xB35C) 、功 (0xA55C) 、蓋 (0xBB5C)

Unicode VS ISO 10646

- ❑ 1991年左右，同時有兩個組織著手規範世界字碼
 - Unicode
 - ISO 10646

- ❑ 過不了多久，他們就互相體認到「這個世界不需要兩套不同的單一字符集」
 - 因此他們決定共用同樣的字碼

- ❑ 現在這兩個組織各自存在，各自互相砥礪

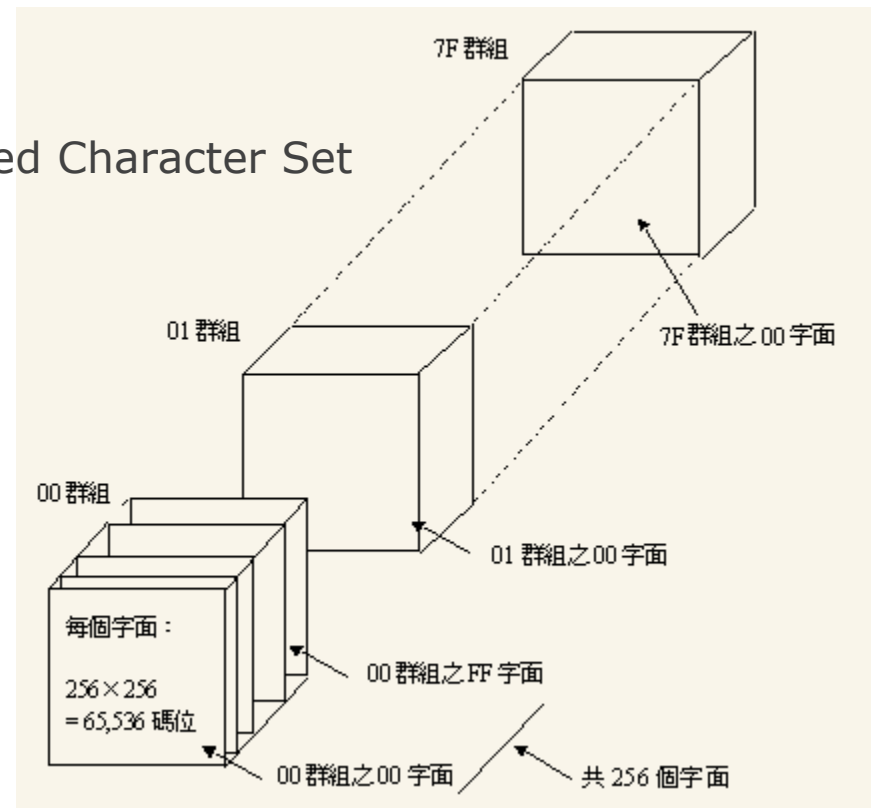
ref:http://zh.wikipedia.org/wiki/ISO_10646

編碼標準 –

ISO10646 and Unicode (1)

□ Goal

- 集結全球通用字符集,成一大聯集
- UCS-4
 - Universal multiple-octet coded Character Set
- 4 bytes encoding (2^{31})
 - 128 Groups
 - 256 Planes each group
 - 256 Rows each plane
 - 256 Cells each row
- BMP (UCS-2)
 - Basic Multilingual Plane
 - 00 group, 00 plane
 - 65536 encoding space
- Why in BMP
 - 若所有字集都在 **BMP** 中, 就可以只使用 **2 bytes**, 否則就要用 **4 bytes**, 不能混用



編碼標準 - ISO10646 and Unicode (2)

□ BMP

列八位元組	基本拉丁文		拉丁文 1 補充	
00	基本拉丁文		拉丁文 1 補充	
01	拉丁文擴充 A		拉丁文擴充 B	
02	拉丁文擴充 B	國際音標擴充	間隔修飾字元	
03	結合之附加記號		基本希臘文	希臘符號和哥普特文
04	斯拉夫文字母			
05	亞美尼亞文		希伯來文 (基本和擴充)	
06	基本阿拉伯文		阿拉伯文擴充	
09	古梵文		孟加拉文	
0A	錫克教文		印度文	
0B	印度文		坦米爾文	
0C	德拉威 Telugu 文		德拉威 Kannade 文	
0D	德拉威 Malayalam 文			
0E	泰文		寮文	
0F			基本藏文	
10			喬治亞文	
11	韓文拼音符號 (Hangul Jamo)			
1E	拉丁文擴充附加			
1F	希臘文擴充			
20	一般標點符號	二/下標	錢幣符號	與符號組合之附加記號
21	似字母的符號		數字形式	箭號
22	數學運算符			
23	其他技術符號			
24	控制圖象	光學字元識別	括號文數字	
25	製表格圖	區塊元件	幾何形狀	
26	其他符號			
27	什錦符號			
30	中日韓符號和標點		平假名	片假名
31	注音符號	韓文相容拼音	中日韓其他字元	
32	中日韓括號字母和月份			
33	中日韓相容字元			
34	中日韓認同的表意文字擴充 A (CJK Unified Ideographs Extension A)			
4D				
4E	中日韓認同的表意文字 (CJK Unified Ideographs)			
9F				
A0				
AB				
AC	韓文拼音(Hangul)			
D7				
D8	(UTF-16 使用區)			
DF				
E0	專用區			
F8				
F9	中日韓相容的表意文字			
FA				
FB	字母表現形式			
FC	阿拉伯文表現形式 A			
FD				
FE	組合半形標示	中日韓相容形式	小寫變體	阿拉伯文表現形式 B
FF	半形和全形		特殊符號	

A 區

I 區

O 區

S 區

R 區

Unicode Transformation Format

□ UTF: UCS/Unicode Transformation Format

- UTF-16(2、4 bytes)
 - 將一個 32-bit ISO10646 字元轉成多個 16-bit Unicode
- UTF-8(1~4 bytes)
 - 將一個32-bit ISO10646 字元轉成多個 8-bit Unicode
 - 將一個16-bit Unicode 字元轉成多個 8-bit Unicode
 - 128個US-ASCII字元只需1 bytes編碼

非常經典的 UTF-8...

- ❑ 與既有系統的相容性
 - 只包含 ASCII 0-127 的字串是合法的 UTF-8 字串
 - NULL-terminated 字串處理

- ❑ 極高的辨識性
 - UTF-8字串可以由一個簡單的演算法可靠地識別出來。

- ❑ 可以容納所有 Unicode 字元
 - UTF-8 理論值可以容納百萬個字元 (實際是 1112064 個)
 - (2012 年發佈的 Unicode 6.2 也才十一萬個字元)

- ❑ Unicode 與 UTF-8 之間的轉換很方便

中文環境 (1)

□ 要做到哪些事情

- 中文訊息
 - 中文顯示
 - 中文輸入
 - 中文列印
 - 中文處理
- 簡單
- ↓
- 困難

中文環境 (2)

□ 中文化方式

- 直接修改程式
 - 套件以排山倒海之勢而來
- **18 chars**
- 國際化(**I**nternationalization**N** , i18n)
 - Multi-language architecture
 - 程式設計人員按照該架構的機制與準則寫程式, 便可支援各式各樣的語言
 - Locale (**LOCAL**ization Environment database)
 - 程式根據使用者選擇的 locale 聯繫到不同資料庫, 進而提供該語言的支援
- 在地化(**L**ocalization**N** , L10n)
 - 在 i18n 的大架構下, 加入「在地化」的特性
- 通常i18n只需做一次, 而L10n要針對每個語言個別做

中文環境 (3)

□ locale in FreeBSD

- 地區性語言的資訊
 - LC_ALL
 - 掌管該 locale 中所有字元的處理方式
 - LC_CTYPE
 - 掌管程式訊息輸出所用的語言
 - LC_MESSAGES
 - 掌管程式訊息輸出所用的語言
 - LC_TIME
 - 時間格式
 - LC_NUMERIC
 - 數字格式
 - LC_MONETARY
 - 貨幣格式
 - LC_COLLATE
 - 字母順序與特殊字元比較
 - LANG
 - 語言顯示
- 效力優先性：LC_ALL > LC_* > LANG

中文環境 (4)

□ 設定 locale

- csh/tcsh shell
 - `setenv LC_CTYPE en_US.UTF-8`
- Bourne Shell
 - `export LC_CTYPE=en_US.UTF-8`
- `/usr/share/locale/`
 - 各國的 locale 資訊
 - 命名規則：語言_地區名.字元編碼名稱
 - `zh_TW.UTF-8`
 - `zh_CN.GBK`

中文環境 (5)

- 中文 Terminal (Remote Login)
 - M\$ Windows: putty, piety, netterm, multi-term, telnet, ...etc.
 - X Window: xterm, rxvt, aterm, mterm, roxterm...etc.
 - 設定好中文支援，登入後即可看到中文

Steps of Exercise

□ 中文 Xwindow

- 安裝中文字型
- 設定 Shell locale 環境
- 安裝中文 Terminal
- 安裝 ibus 中文輸入程式

安裝中文字型 (1)

□ 兩大中文字型種類

- 點陣字型 (Bitmapped Font)
 - BDF (Bitmap Distribution Format) 點陣分散格式
 - HBF (Hanzi Bitmap Font) 漢字點陣字體
 - PCF (Portable Compiled Font)
- 曲線描邊字型 (Outline Fonts)
 - True Type Font (TTF)

安裝中文字型 (2)

□ Font Path

Font Path:

```
/usr/local/lib/X11/fonts/misc/  
/usr/local/lib/X11/fonts/TTF/  
/usr/local/lib/X11/fonts/Type1/  
/usr/local/lib/X11/fonts/75dpi/  
/usr/local/lib/X11/fonts/100dpi/  
/usr/local/lib/X11/fonts/local/
```

□ 安裝字型

- 透過 **ports** 安裝字型檔案
- 使用 **ttfm** 安裝該字型
- 使用 **fc-cache** 建立字型資料庫
- 修改各軟體設定使用別的字型

安裝中文字型 (3)

- ❑ 安裝 ttfm – TrueType 字型管理工具
- ❑ ttfm
 - ttfinfo 讀取 ttf 字型格式資訊的程式
 - **% ttfinfo /usr/local/share/fonts/TrueType/fireflysung.ttf**

```
(21:38)wjguo@[oopc6:/home/wjguo] >ttfinfo /usr/local/share/fonts/TrueType/fireflysung.ttf
TTFINFO_FONT_FILE="/usr/local/share/fonts/TrueType/fireflysung.ttf"
TTFINFO_FACE_NUM="1"
TTFINFO_FACE_INDEX="0"
TTFINFO_FONT_NAME="AR PL New Sung"
TTFINFO_FONT_PSNAME="AR-PL-New-Sung"
TTFINFO_FOUNDRY_NAME="misc"
TTFINFO_WEIGHT_NAME="medium"
TTFINFO_WIDTH="normal"
TTFINFO_NUMCMAP="3"
TTFINFO_CMAP0="0,3"
TTFINFO_CMAPNAME0="Apple Unicode,(v.2.0)"
TTFINFO_CMAP1="1,0"
TTFINFO_CMAPNAME1="Apple,Roman"
TTFINFO_CMAP2="3,1"
TTFINFO_CMAPNAME2="Windows,Unicode"
TTFINFO_MAPNUM="3"
TTFINFO_FONTMAP1="-misc-AR PL New Sung-medium-r-normal--0-0-0-0-p-0-big5-0"
TTFINFO_FONTMAP2="-misc-AR PL New Sung-medium-r-normal--0-0-0-0-p-0-gb2312.1980-0"
TTFINFO_FONTMAP3="-misc-AR PL New Sung-medium-r-normal--0-0-0-0-p-0-jisx0208.1983-0"
```

安裝中文字型 (4)

- **ttfm.sh**

```
(21:38)wjguo@[oopc6:/home/wjguo] >ttfm.sh
```

```
True-Type Font Manager 0.9.3
```

```
Usage: /usr/local/bin/ttfm.sh [option]
```

```
--add [module] <file>...    install ttf font
--remove [module] <file>...  remove ttf font from the system
--list <module>...          list all ttf fonts on the system
--modules                   list all ttf manager modules on the system
--setdefault <module> <file>
                             set default ming font of module to file
--setdefault_kai <module> <file>
                             set default kai font of module to file
--initm <module>..          initialize modules
--help                       show this info
```

安裝中文字型 (5)

□ 選一個來裝 ...

Chapter 6. 輸出字型

Table of Contents

- 6.1. [Bitmapped Font - 點陣字型概論](#)
- 6.2. [cmexfonts - 中推會 Big5+ 點陣字型](#)
- 6.3. [kcfnts - 國喬點陣字型](#)
- 6.4. [gugod-clean - 搭配中文點陣字型用的英文點陣字型](#)
- 6.5. [intlfonts - 各國的免費點陣字型](#)
- 6.6. [PostScript 概論](#)
- 6.7. [使用 TrueType 字型當作是 CID fonts](#)
- 6.8. [moefonts-cid - 由 Adobe 轉譯的 MOE CID Font](#)
- 6.9. [以 gs 觀看不內嵌的 pdf 檔](#)
- 6.10. [TrueType - 全真字型概論](#)
- 6.11. [tffm - TrueType 字型管理工具](#)
- 6.12. [mingliu - 微軟細明體 TrueType 字型](#)
- 6.13. [simsun - 微軟宋體 TrueType 字型](#)
- 6.14. [mingunittf - 香港補增字符集2001](#)
- 6.15. [arnettf](#)
- 6.16. [fireflytff - 內嵌點陣字的自由字型](#)
- 6.17. [moettf - 台灣教育部標準 TrueType 字型](#)
- 6.18. [arphicttf - 文鼎科技 TrueType 字型](#)
- 6.19. [wangtff - 王漢宗教授 TrueType 字型](#)
- 6.20. [ntuttff - 台大字型](#)
- 6.21. [oto - Open Type Organizer 程式](#)

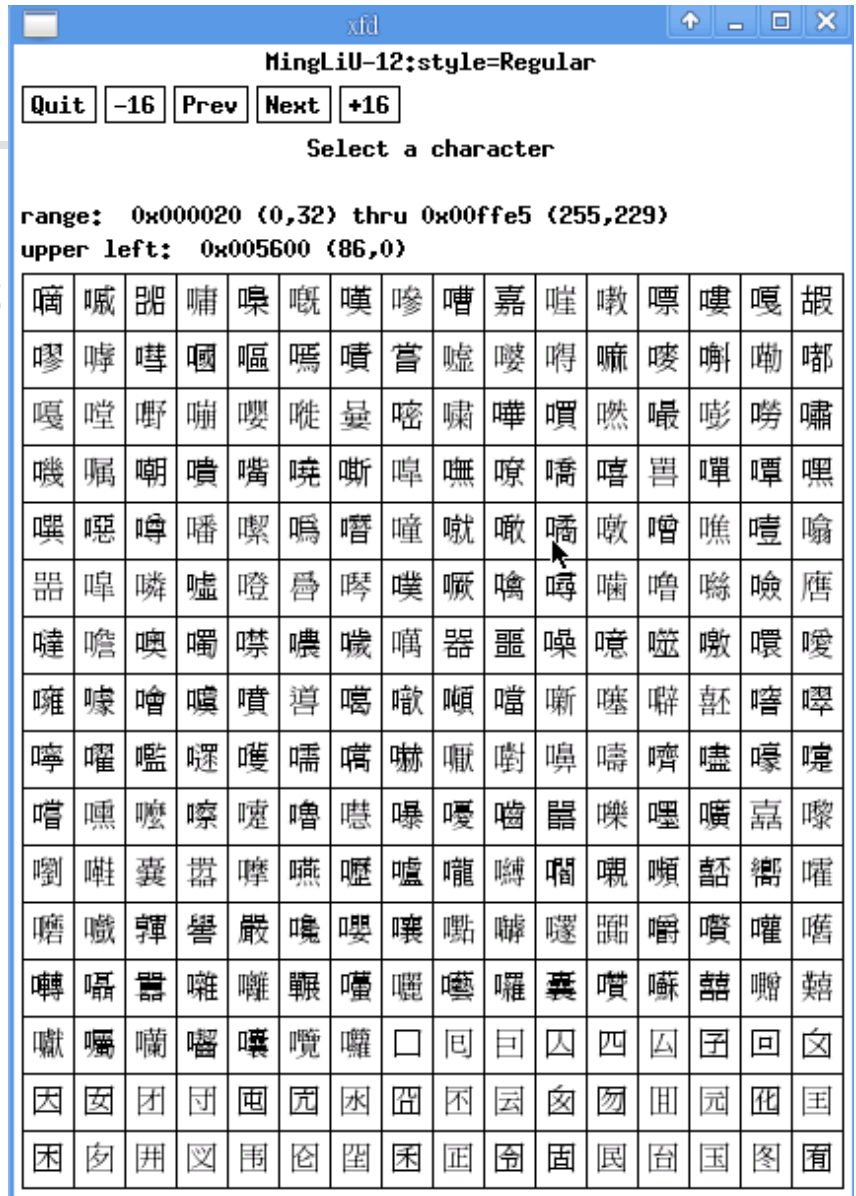
安裝中文字型 (6)

□ Fireflyttf

- 透過 `ports` 安裝的都會自己跑
 - `% ttfm.sh --add xttfm /usr/local/share/fonts/TrueType/fireflysung.ttf`
 - `% fc-cache -f -v /usr/local/lib/X11/fonts/TrueType/`
- `portmaster chinese/fireflyttf`

*安裝中文字型 (7)

- 用 xfd 來看
 - % setenv LC_CTYPE zh_TW.Big5
 - % xfd -fa "MingLiU"
 - man xfd



安裝中文字型 (8)

□ 增加 Font Path

- Edit /etc/X11/xorg.conf

- /usr/local/share/fonts/TrueType/fireflysung.ttf

- /usr/local/lib/X11/fonts/TrueType/fireflysung.ttf

symbolic link



- Restart xwindow

Section "Files"

ModulePath "/usr/local/lib/xorg/modules"

FontPath "/usr/local/lib/X11/fonts/misc/"

FontPath "/usr/local/lib/X11/fonts/TTF/"

FontPath "/usr/local/lib/X11/fonts/OTF"

FontPath "/usr/local/lib/X11/fonts/Type1/"

FontPath "/usr/local/lib/X11/fonts/100dpi/"

FontPath "/usr/local/lib/X11/fonts/75dpi/"

FontPath "/usr/local/lib/X11/fonts/TrueType/"

FontPath "/usr/local/lib/X11/fonts/local/"

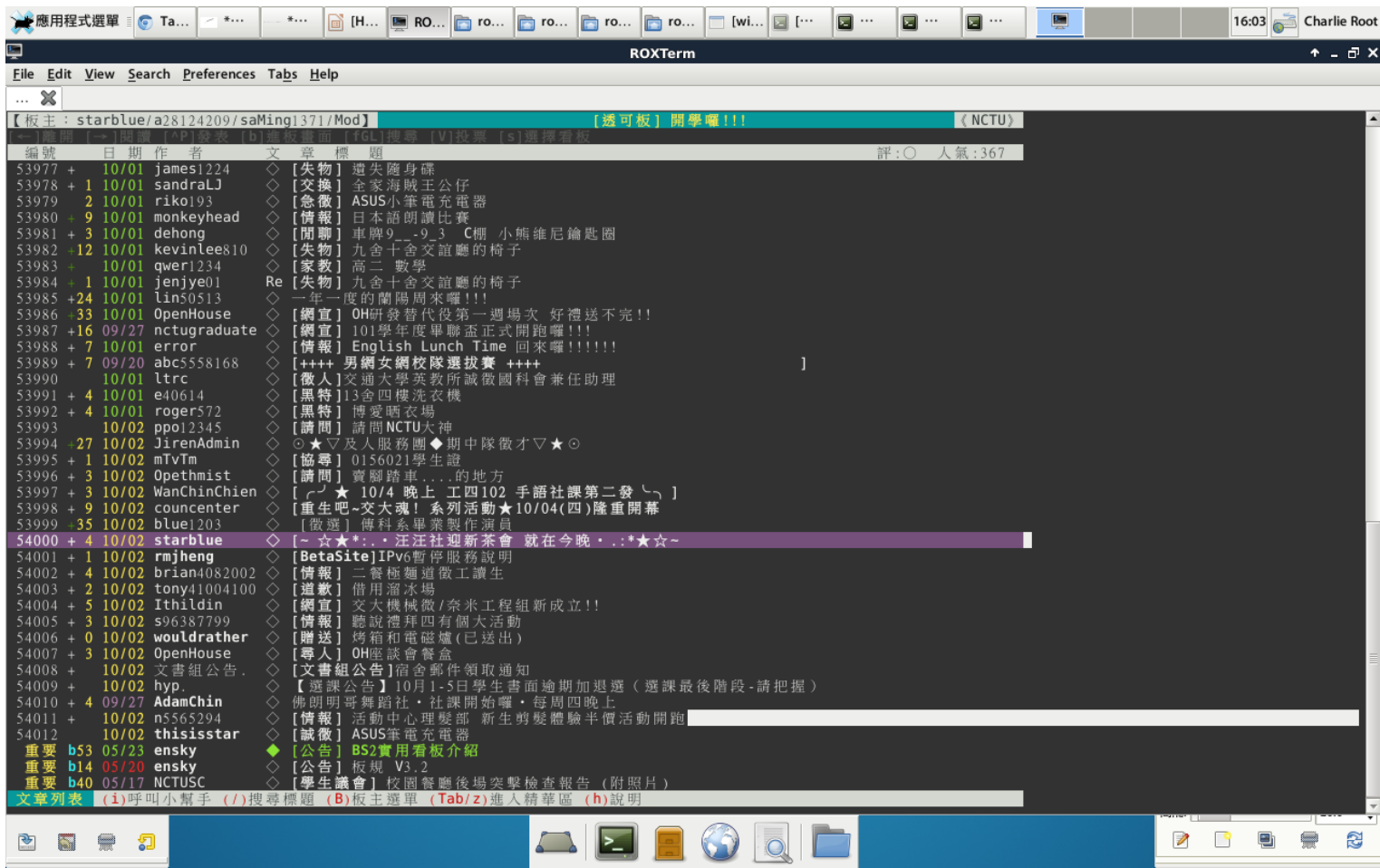
EndSection

安裝中文 Terminal

- ❑ rxvt
 - /usr/ports/x11/rxvt-unicode
- ❑ aterm
 - /usr/ports/chinese/aterm
- ❑ eterm
 - /usr/ports/chinese/eterm
- ❑ ROXterm
 - /usr/ports/x11/roxterm
- ❑ mlterm
 - /usr/ports/x11/mlterm

ROXterm

- ❑ X11/roxterm
- ❑ roxterm-config



安裝中文輸入程式

- ❑ Choicesibus-chewing(chinese/ibus-chewing)
 - ibus-pinyin(chinese/ibus-pinyin)

安裝 ibus 中文輸入程式 (1)

❑ ibus

- Intelligent Input Bus

1. `setenv LC_CTYPE zh_TW.UTF-8` (csh/tcsh)
`export LC_CTYPE=zh_TW.UTF-8` (sh/bash)2.
2. Edit `.xinitrc`

```
XIM=ibus
GTK_IM_MODULE=ibus
QT_IM_MODULE=xim
XMODIFIERS=@im=ibus'
XIM_PROGRAM="ibus-daemon"
XIM_ARGS="--daemonize --xim"
```

安裝 ibus 中文輸入程式 (2)

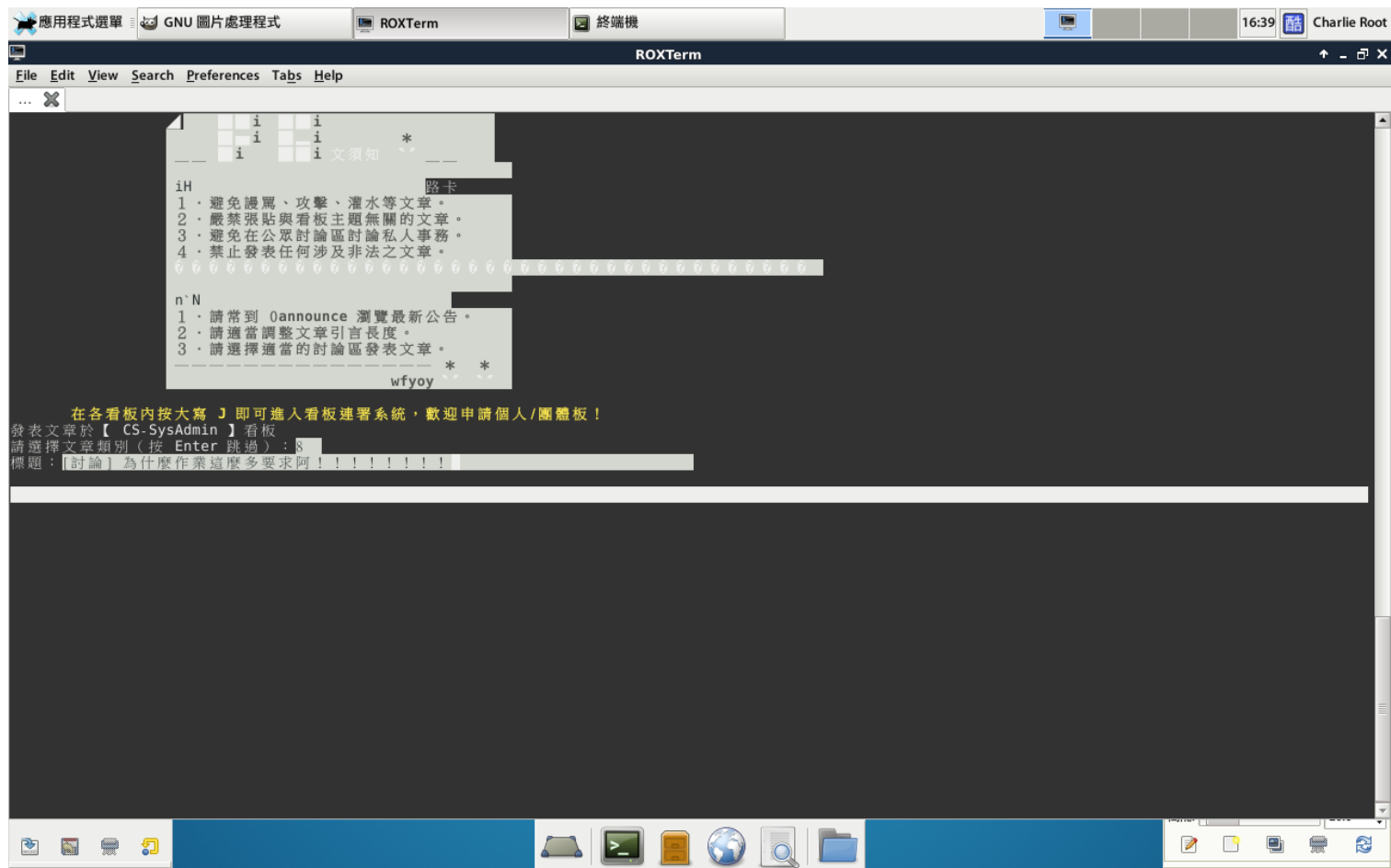
❑ ibus相關設定

- % ibus-setup
- 加入 Chewing



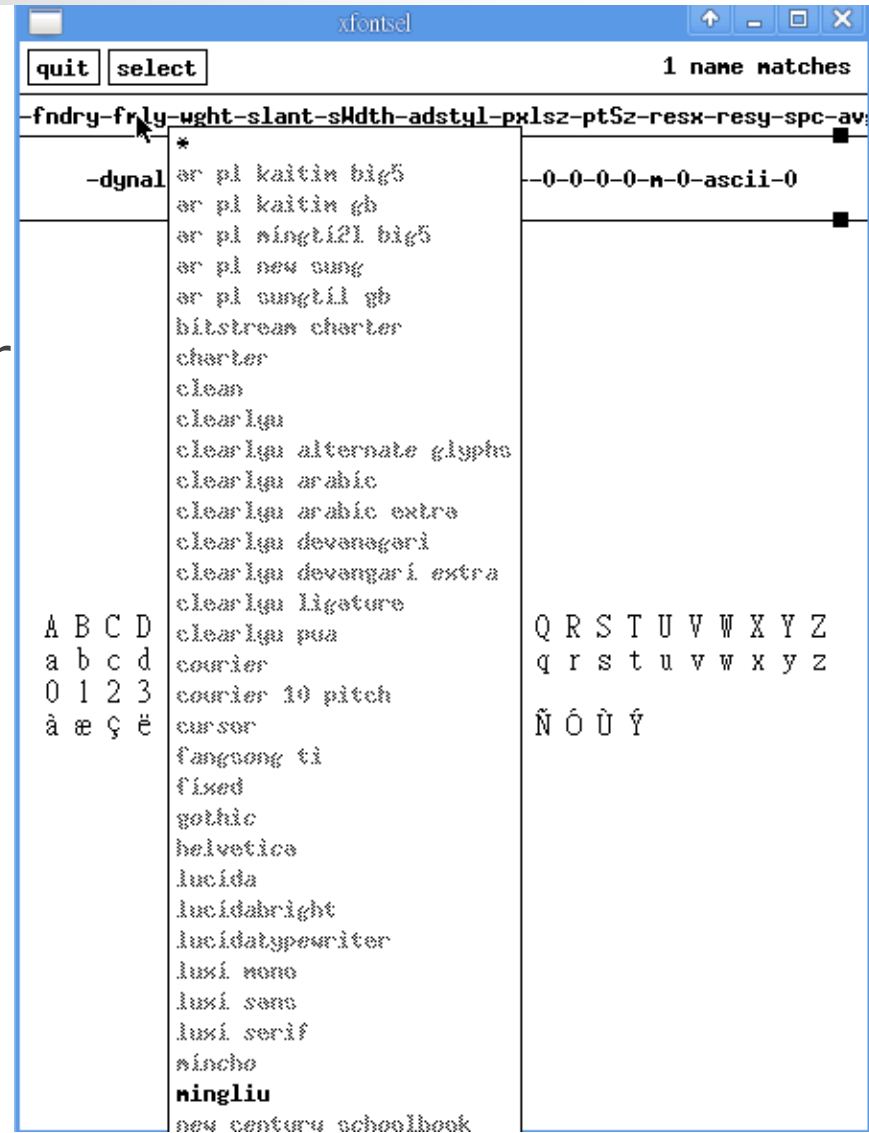
安裝 ibus 中文輸入程式 (3)

4. Switch to chinese input: Ctrl -Space



其他設定 (1)

- 顯示所有可用字型
 - xlsfont
- 選擇字型程式
 - xfontsel – X font selector
 - click [select]
 - 滑鼠中間鍵貼上
 - 滑鼠左右鍵貼上



References

☐ 中文碼介紹

- <http://www.cns11643.gov.tw/web/word.jsp>

☐ **FreeBSD Chinese HOWTO**

- <http://netlab.cse.yzu.edu.tw/~statue/freebsd/zh-tut/index.html>

☐ **Introduction to i18n**

- <http://www.debian.org/doc/manuals/intro-i18n/>