

MATBN 中文廣播新聞語料庫 簡介

民國 90 年至 93 年間，國內從事語音研究之相關學校及研究單位聯合向國科會中提出一語料蒐集計畫－『中文自發性語音語料庫之建立』(Spontaneous Mandarin Speech: Corpus and Processing；計畫編號：NSC-92-2213-E-009-021)，參與的單位共有國立交通大學電信工程學系、國立台灣大學電機工程學系、國立清華大學電機工程學系、國立成功大學電機工程學系、中央研究院資訊研究所、工研院前瞻研究中心及中華電信研究所。在該項『中文自發性語音語料庫之建立』計畫中收錄了新聞語音資料，語料來源由公共電視台之新聞，計畫中並對所蒐集之 197 個小時節目之語音資料做人工文字標註(轉記；transcription)處理。會讓該項成果可與國內外從事國語語音研究之單位分享，今將申請將上述之研究成果申請技轉中華民國計算語言學學會(The Association for Computational Linguistics and Chinese Language Processing；ACLCLP)。

語音資料庫相關資訊

1. 授權語音資料庫共包含 197 個小時之語音資訊及其內容之轉記與語音資訊之標註資料；共五片 DVD 光碟。

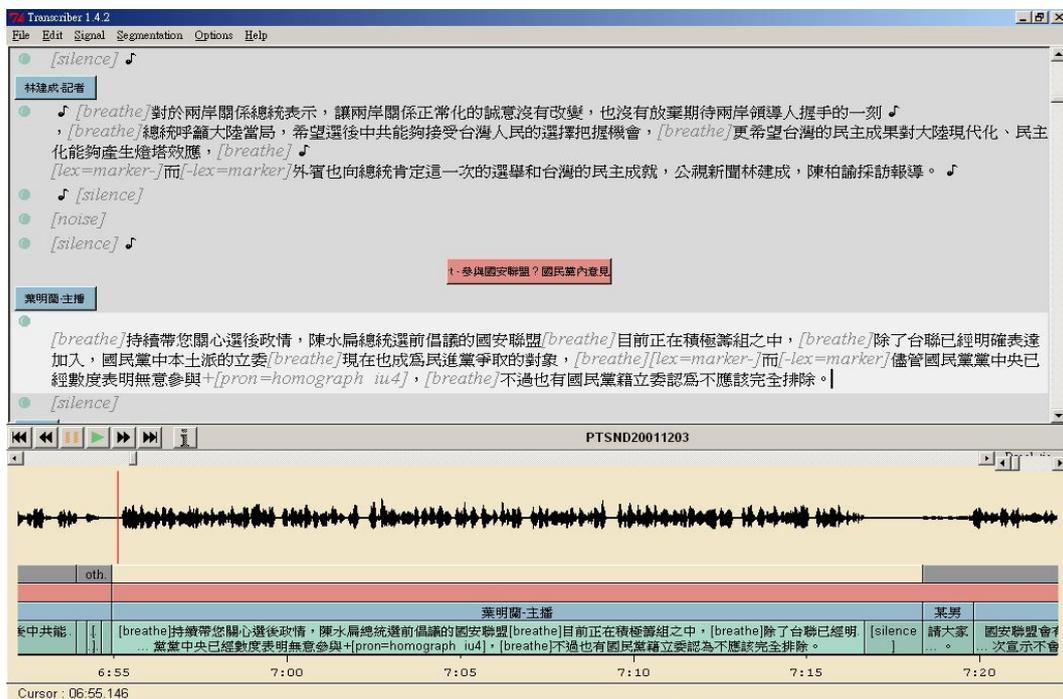
2. 錄製語料的相關資訊：

所有的新聞語音都在公視錄製的，在第一年度計畫執行之初的數個月先進行各項準備工作，包括聯繫電視及廣播公司洽談授權、準備標註軟體、決定標註方式等。經與公共電視洽談後，公視同意授權我們使用其新聞節目，並建議我們採用『公視新聞深度報導』節目及願意協助我們錄音(影)，錄音工作自 90 年 11 月 7 日起正式展開。『公視新聞深度報導』於每週一至五晚間 21:00-22:00 播出一個小時，自 91 年 7 月起，變更節目名稱為『公視晚間新聞』，自 91 年 9 月起，播出時間改為晚間 21:00-21:45，播出 45 分鐘，另於 21:45-22:00 播出 15 分鐘的『公視手語新聞』，92 年 1 月 31 日起，『公視晚間新聞』移至 19:00-20:00 播出，21:00-21:45 則播出『公視全球現場』，21:45-22:00 仍播出『公視手語新聞』。自 90 年 11 月 7 日起至 92 年 2 月底止，錄音時間固定為 21:00-22:00，92 年 3 月起，錄音時間則包括 19:00-20:00 及 21:00-22:00 兩個時段。本計畫錄音工作進行至 92 年 6 月底結束，共

收錄約 300 個小時的新聞節目，主要內容為國內新聞，也有一小部分為國際新聞。語料已轉為 windows PCM 的聲音檔，其規格為為 16kHz、16 bit、單軌聲音檔。

3. 語料資料之人工文字標註資料

語音的標註是由受過訓練的專任助理進行較準確的文字標註，並做 cross checking，此部分工作由中央研究院資訊所王新民博士統籌規劃及負責推動。我們採用 LDC (Linguistic Data Consortium) 提供的 Transcriber 系統來標註電視新聞錄音資料，請參考圖一。首先，將公視取回的 DAT(數位錄音帶)，經 USB 介面直接將錄音帶內的數位信號讀進 PC 內轉為格式為 44.1kHz、16bit、stereo 的聲音檔 (windows PCM、.wav)，並燒錄於光碟中以便保存。然後，將檔案轉成標註使用的聲音檔，因考量檔案傳輸及讀取速度的問題，將原始的檔案，利用聲音編輯軟體 — CoolEdit 2000 將已轉為 windows PCM 的聲音檔進行格式轉換。轉換為 16kHz、16 bit、mono 後，為便利日後管理及利用，每週的公視新聞深度報導，每月的公視演講廳、客家新聞雜誌分別儲存於同一光碟中保存。



圖一：利用 Transcriber 標註新聞語音的實例

在標註過程中，舉凡雜訊、背景環境、發音不標準、方言、說話者性別、主播/記者/被採訪者等資訊都盡量鉅細靡遺標註下來，標註的結果以 XML 檔案儲存，請參考圖二。

```

UltraEdit-32 - [E:\whm\project\speechcollection\data\PTSND20011203\PTSND20011203.TRS]
檔案(F) 編輯(E) 搜尋(S) 專案(O) 檢視(V) 格式(I) 行列(L) 巨集(M) 進階(A) 視窗(W) 輔助(H)
PTSND20011203.TRS
外賓也向總統肯定這一次的選舉和台灣的民主成就，公視新聞林建成，陳柏諭採訪報導。
<Background time="413.931" type="other" level="off"/>
<Sync time="413.932"/>
<Background time="413.932" type="other" level="high"/>
<Event desc="silence" type="noise" extent="instantaneous"/>
<Sync time="414.239"/>
<Event desc="noise" type="noise" extent="instantaneous"/>
<Sync time="414.642"/>
<Event desc="silence" type="noise" extent="instantaneous"/>
<Background time="414.967" type="other" level="off"/>
</Turn>
</Section>
<Section type="report" topic="to8" startTime="414.967" endTime="544.297">
<Turn speaker="spk1" mode="planned" fidelity="high" channel="studio" startTime="414.967" endTime="438.353">
<Sync time="414.967"/>
<Event desc="breathe" type="noise" extent="instantaneous"/>
持續帶您關心選後政情，陳水扁總統選前倡議的國安聯盟
<Event desc="breathe" type="noise" extent="instantaneous"/>
目前正在積極籌組之中，
<Event desc="breathe" type="noise" extent="instantaneous"/>
除了台聯已經明確表達加入，國民黨中本土派的立委
<Event desc="breathe" type="noise" extent="instantaneous"/>
現在也成為民進黨爭取的對象，
<Event desc="breathe" type="noise" extent="instantaneous"/>
<Event desc="marker" type="lexical" extent="begin"/>

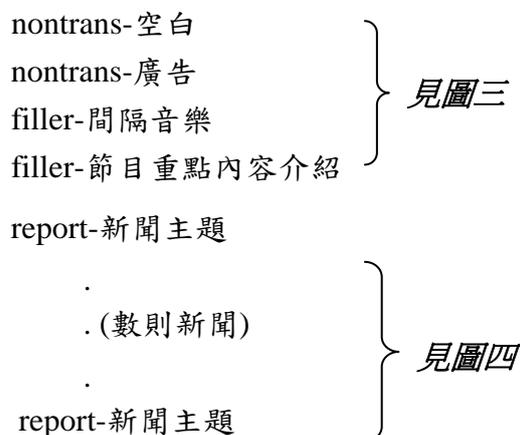
```

圖二：Transcriber 的 XML 標註檔案

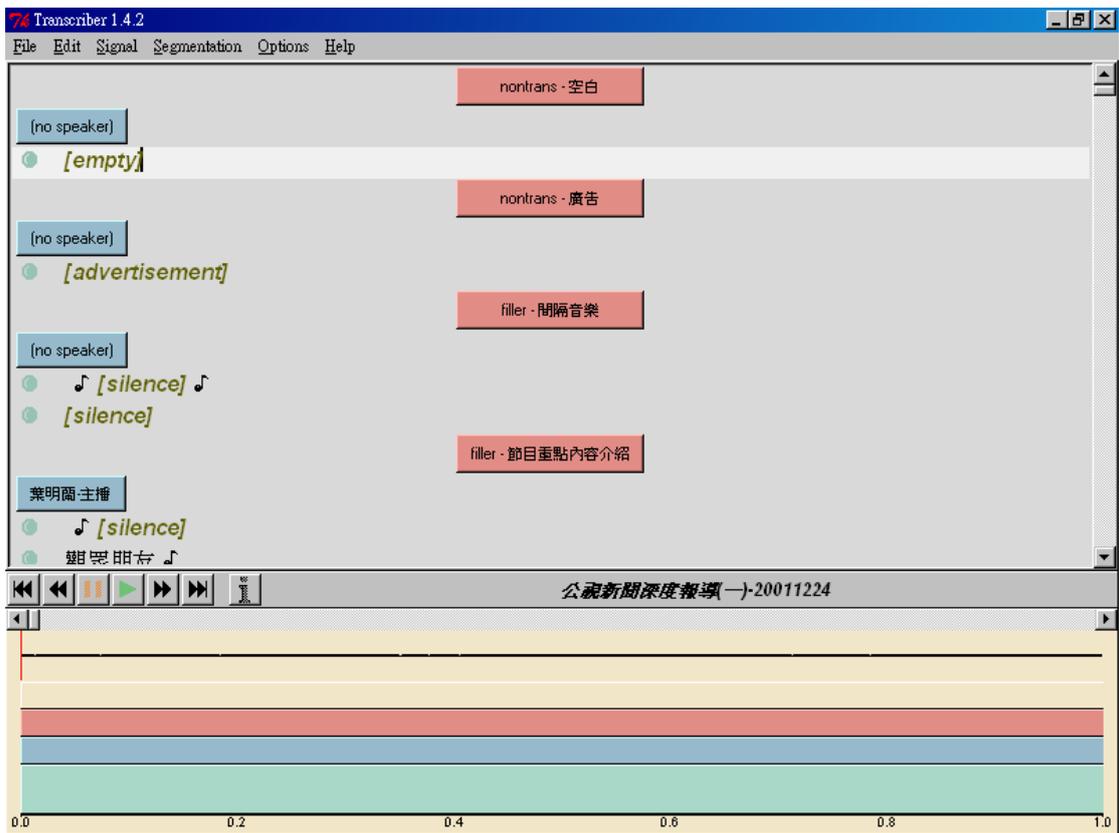
標註重點主要分為四大部分，分別為：

- ◆ 段落主題
- ◆ 說話者名稱
- ◆ 背景聲音
- ◆ 插入事件

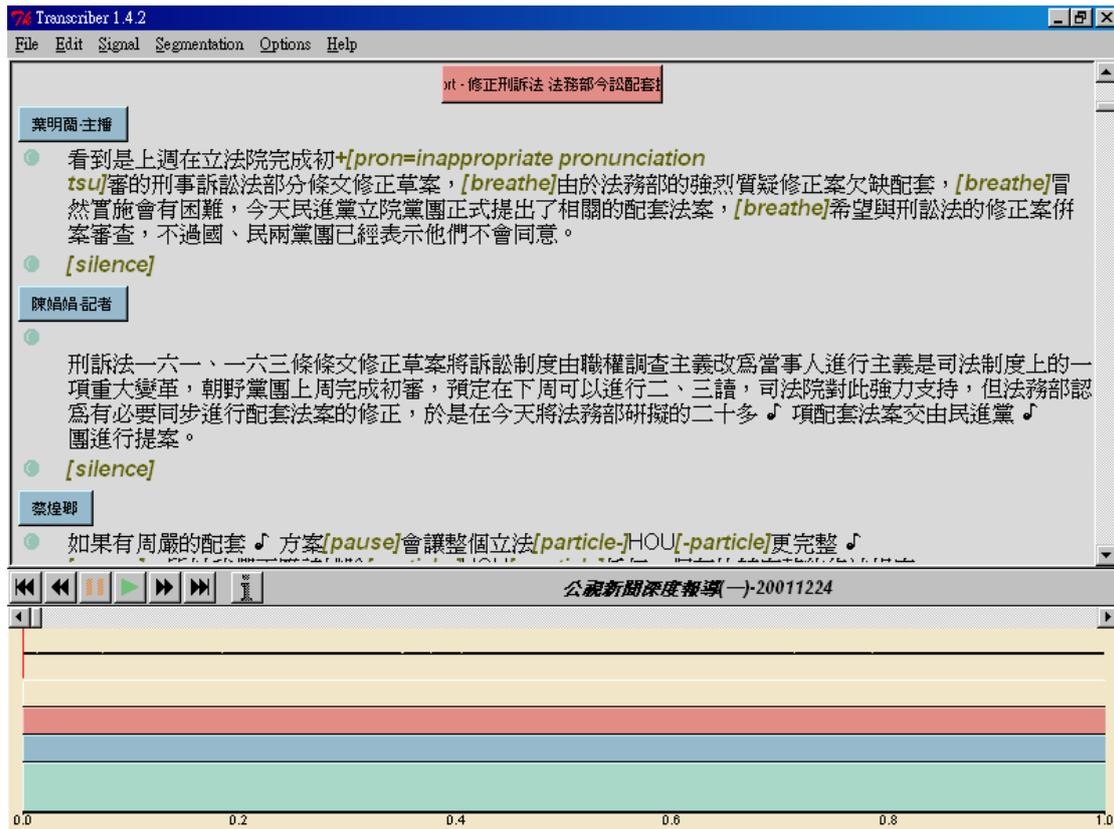
『公視新聞深度報導』於 Transcriber 系統大致上之基本架構包含：



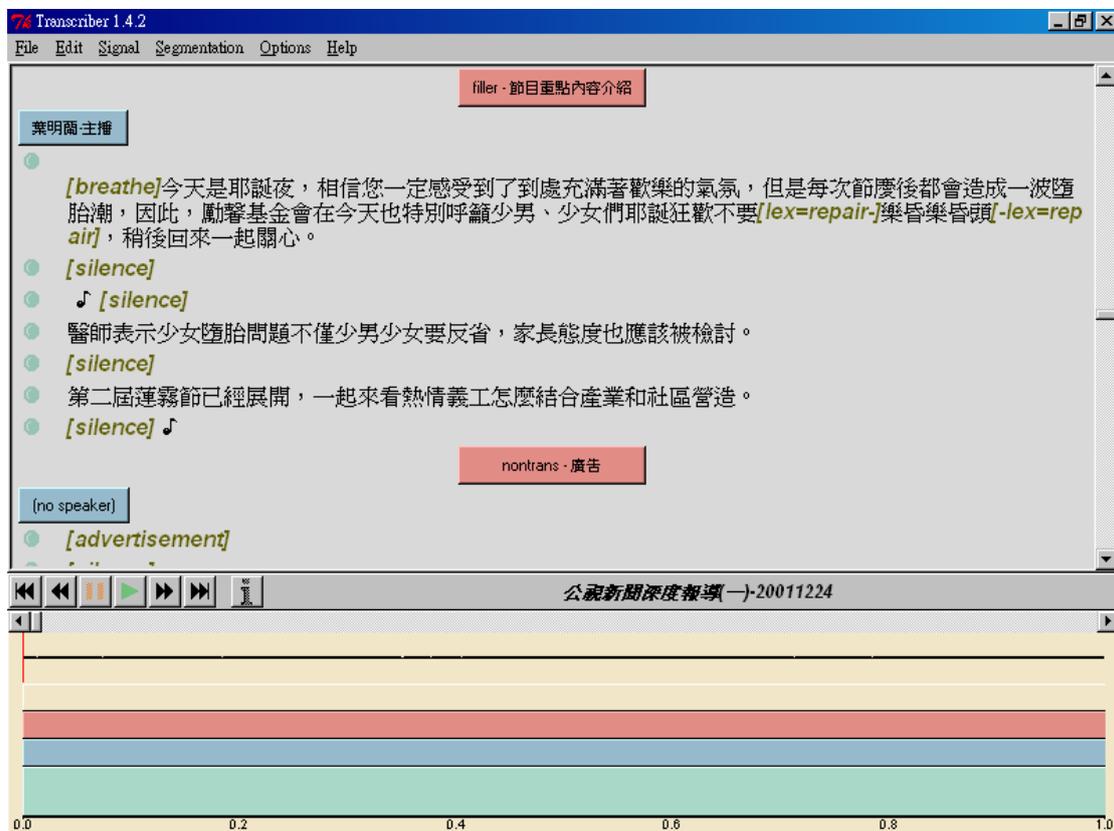
- filler-節目重點內容介紹 } 見圖五
- nontrans-廣告 } 見圖五
- report-新聞主題
- . (數則新聞)
- . report-新聞主題 } 同圖四
- report-氣象預報 } 見圖六
- filler-結尾 } 見圖六
- filler-片尾音樂
- nontrans-廣告 } 見圖七
- nontrans-空白 } 見圖七



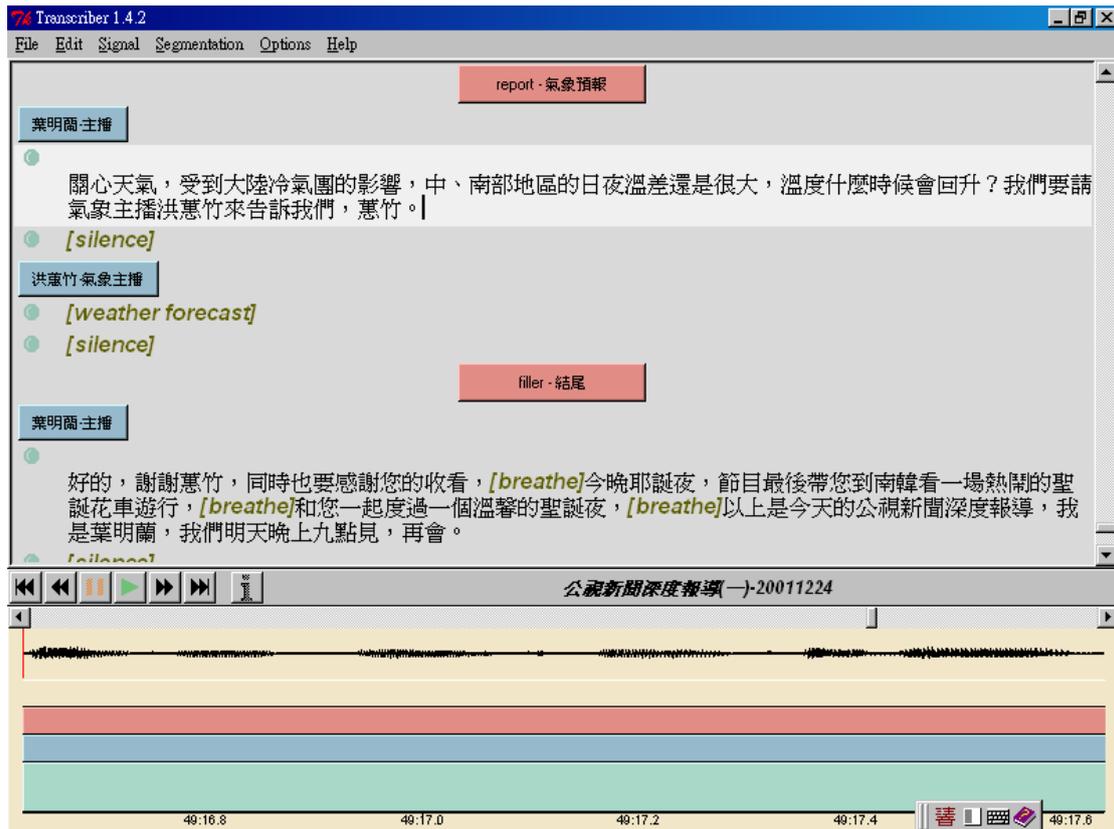
圖三



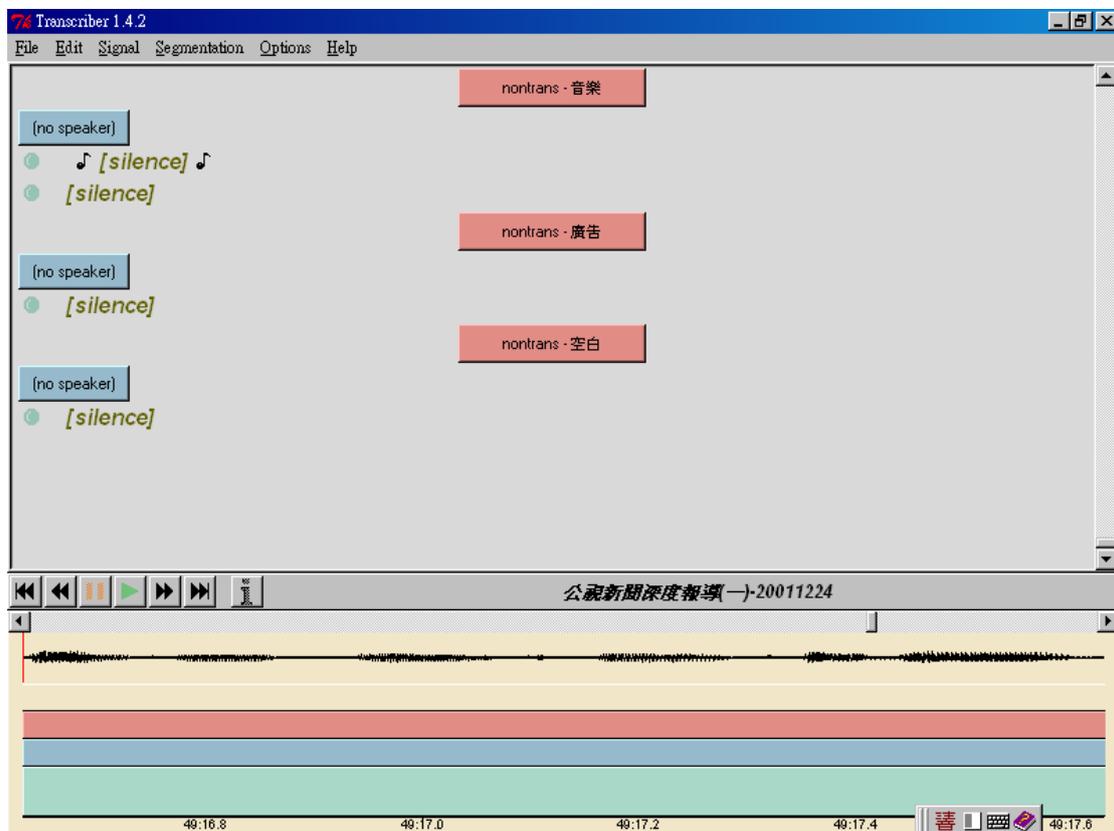
圖四



圖五



圖六



圖七