

Langues peu dotées

Définition et problématiques pour le TALN et le TALP

Vincent Berment

C&S - LIG/GETALP

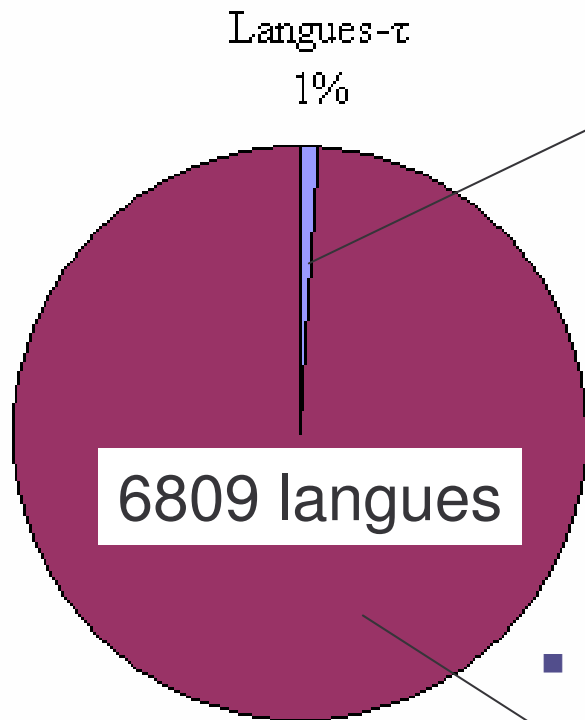
Une mesure du niveau d'informatisation : l'indice- σ

	Services / ressources	Criticité (de 0 à 10)	Note (/20)	Note pondérée (Criticité x Note)
Traitement du texte				
	Saisie simple	10	16	160
	Visualisation / impression	10	14	140
	Recherche et remplacement	8	12	96
	Sélection du texte	6	12	72
	Tri lexicographique	5	0	0
	Correction orthographique	2	0	0
	Correction grammaticale	0	0	0
	Correction stylistique	0	0	0
Traitement de l'oral				
	Synthèse vocale	5	0	0
	Reconnaissance de la parole	5	0	0
Traduction				
	Traduction automatisée	8	4	32
ROC				
	Reconnaissance optique de caractères	9	0	0
Ressources				
	Dictionnaire bilingue	10	4	40
	Dictionnaire d'usage	10	0	0
Total		88		540
Moyenne (/20)				540 / 88 = 6,14

Langues-π et Nations-Unies (Unesco)

- Déclaration universelle sur la diversité culturelle
 - Novembre 2001
 - Extrait : « Toute personne doit ainsi pouvoir **s'exprimer, créer et diffuser ses œuvres** dans la langue de son choix et **en particulier dans sa langue maternelle** »
- Recommandation sur la promotion et l'usage du multilinguisme
 - Octobre 2003
 - Extrait : « Les États membres, les organisations internationales et les entreprises spécialisées dans les technologies de l'information et de la communication devraient appuyer les efforts internationaux de coopération relatifs aux services de **traduction automatisée accessibles à tous**, ainsi qu'aux **systèmes linguistiques intelligents** tels que ceux qui remplissent des fonctions multilingues de **recherche de l'information**, de **dépouillement / résumé** et de **reconnaissance de la parole**, tout en respectant pleinement le droit de traduction des auteurs. »

Langues bien et mal dotées informatiquement



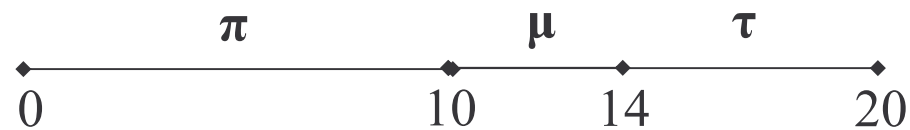
■ Langues Très bien dotées informatiquement

■ Quelques dizaines de langues- τ

- Allemand, anglais, français, japonais, russe...

■ Informatisation rentable => éditeurs de logiciels

- Apple, IBM, Microsoft, Xerox...



■ Langues Peu ou Moyennement dotées

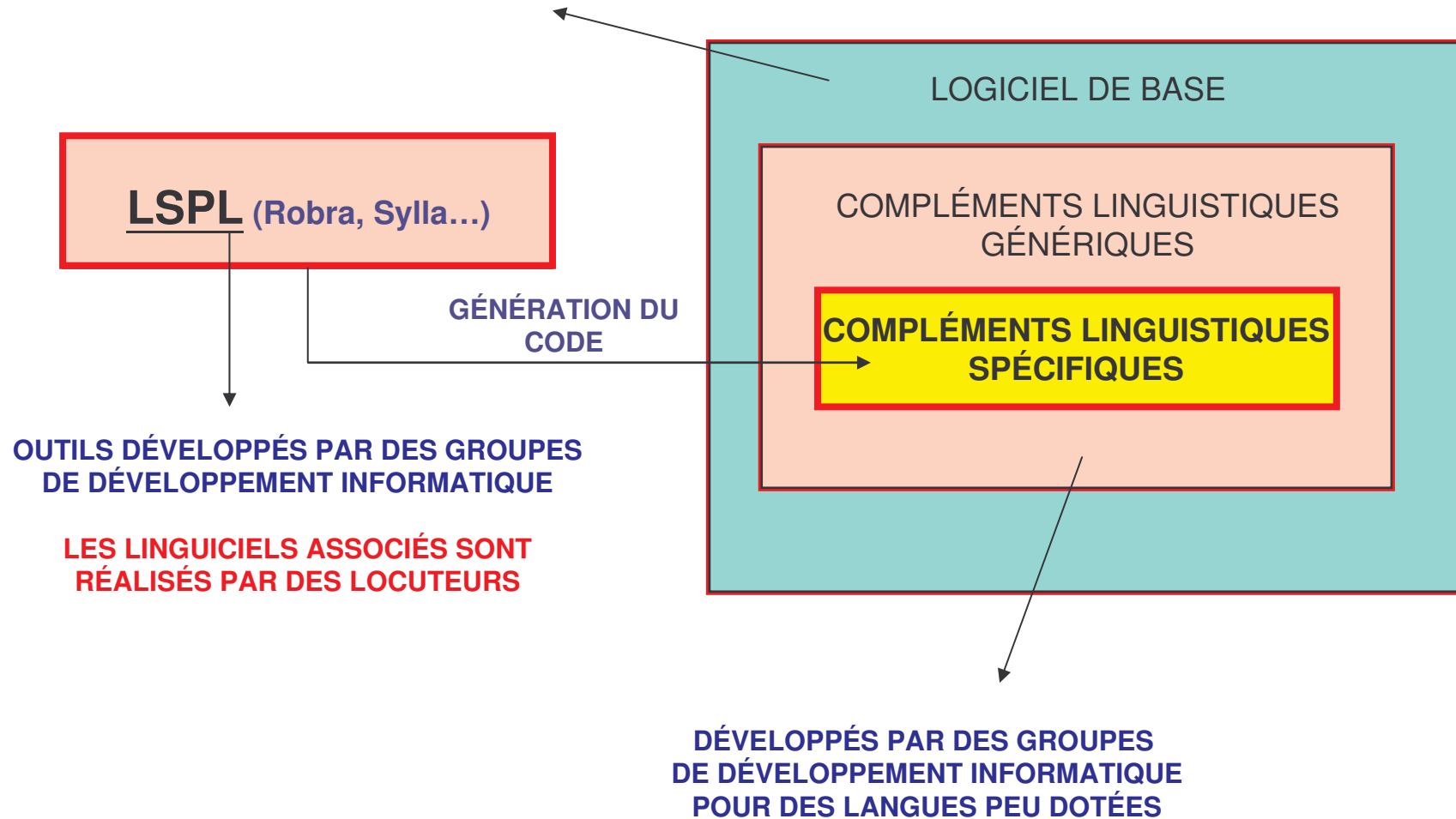
■ Plus de 6000 langues- π et μ

■ Informatisation pas ou peu rentable => autres

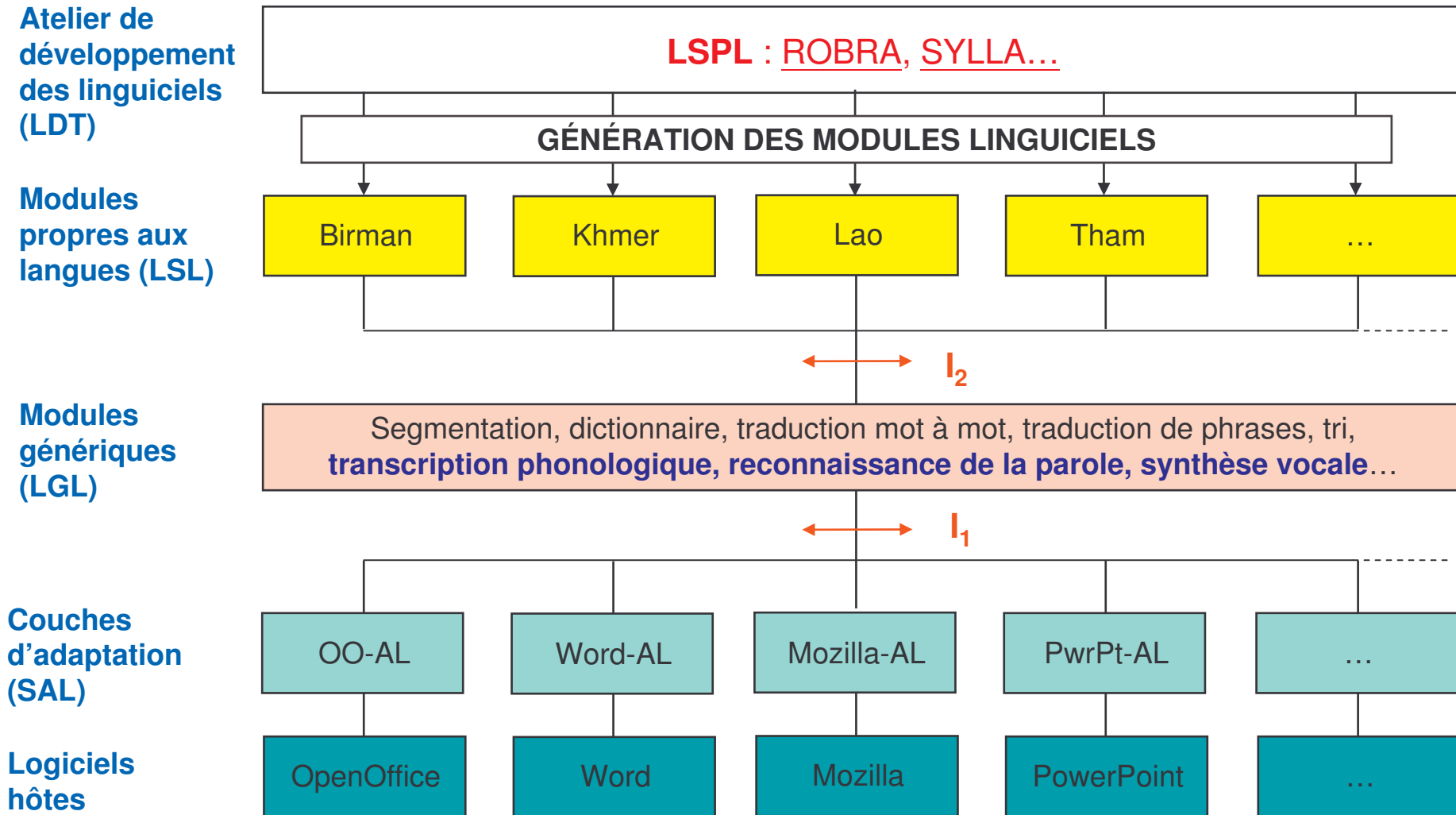
- Groupes de locuteurs créés spontanément
- Projets de développement...

Une méthodologie adaptée

**DÉVELOPPÉ PAR LES GRANDS ÉDITEURS DE LOGICIELS
POUR DES LANGUES BIEN DOTÉES
DOIT PRÉSENTER DES INTERFACES STANDARDISÉES**



Des groupes locaux s'investissent pour leur langue



Gain obtenus pour la saisie grâce à la généricité

Coût la 1 ^{ère} fois (laotien : LaoWord)	Coût les fois suivantes (bengali : BanglaWord)	Gain
250 h	18 h	92,8 %
5 millions de locuteurs au Laos	300 millions de locuteurs au Bangladesh	

Ici, on a substitué le module spécifique au lao par un module spécifique au bengali.

Formule du gain de réutilisation

$$\text{Économie} = (nL * tL * rL * (nE - 1) + nE * tE * rE * (nL - 1)) / (nL * nE * (tL + tE))$$

nL / nE = nombre de langues / d'environnements

tL / tE = temps de développement du code linguistique / générique

rL / rE = taux de réutilisation du code linguistique / générique

- Si $nE=4$, $tL=1000$ heures, $tE=500$ heures, $rL=rE=95\%$

Nb langues	Économie	Économie (en heures)
2	63,33 %	7 600 (4 400 au lieu de 12 000)
5	72,83 %	21 850 (8 150 au lieu de 30 000)
10	76,00 %	45 600 (14 400 au lieu de 60 000)
100	78,85 %	473 100 (126 900 au lieu de 600 000)
1000	79,14 %	4 748 100 (1 251 900 au lieu de 6 000 000)

Un projet d'informatisation de plusieurs langues



- **Projet PAN Localization**
 - Lancé par le Centre de Recherche pour le Développement International (Canada)
 - Mené par le Centre de traitement de l'Ourdou (Pakistan)
- **Budget : 1,3 M€**
- **Informatisation de six langues : bengali, dzongkha (Bhoutan), khmer, lao, népali, cingalais**

Résultats du projet PAN Localization

■ Quelques factorisations :



- Saisie (bengali, lao)
- Conversion Unicode (khmer, cingalais)
- Correcteur d'orthographe (bengali, khmer)
- Tri (bengali , khmer)
- Reconnaissance de caractères (bengali, cingalais)

■ Réussite ?

- + Sensibilisation des équipes locales
- - Non-réutilisation de l'existant venant d'Occident
- - Refus de faire participer les spécialistes occidentaux
- - Résultat de niveau inférieur à l'état de l'art

Autre projet (mené au LIG / GÉTALP / GÉTA)

Segmentation de textes en mots

◦ ນິທານ ເລື່ອງ ນີ້ ກໍ່ ຄວາມ ສະ ເຫຼືອ ນ ໃຈ ແກ່ ຜູ້
ຄົນ ແຕ່ ບູຮານ ນະ ການ ມາ ແລ້ວ. ເລື່ອງ ເລີ່ມ
ຕົ້ນ ຂຶ້ນ ເມື່ອ ແມ່ ຍິງ ສອງ ຄົນ ພາ ກັນ ໄປ ອາບ
ນ້ຳ ຢູ່ ແຄມ ທ່າ, ນ້ຳ ຫ້ວຍ ໃສ່ ດຳ

Ce texte en lao ne contient pas d'espaces (non segmenté)

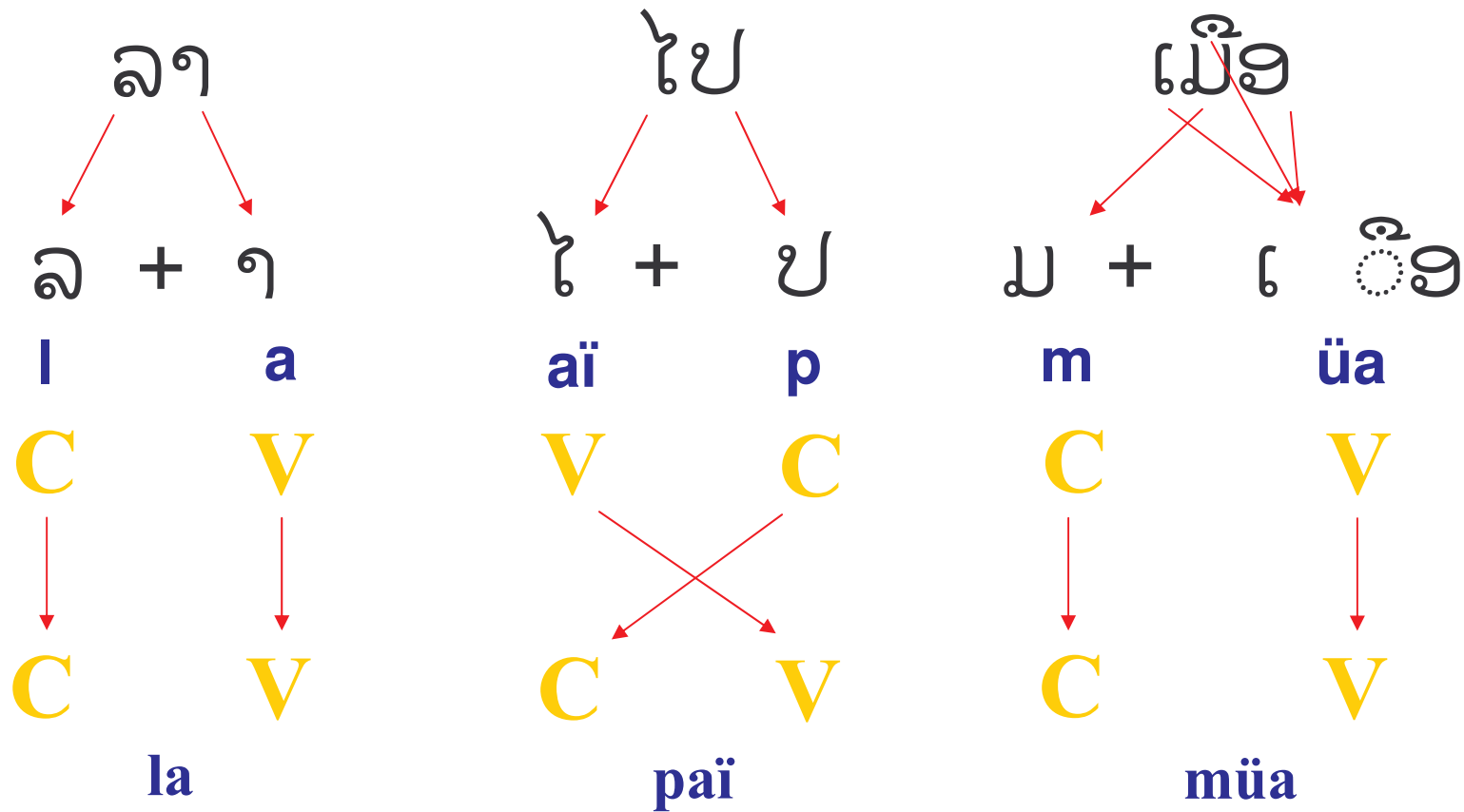
Algorithmes de segmentation en vue de la traduction

- On segmente le texte en syllabes
 - ສະບາຍດີທຸກໆທ່ານ → ສະ-ບາຍ-ດີ-ທຸກ-ໆ-ທ່ານ
- On regroupe les syllabes pour former des mots contenus dans le dictionnaire (algorithme de « plus longue chaîne d'abord »)
 - ສະ-ບາຍ-ດີ-ທຸກ-ໆ-ທ່ານ → ສະບາຍດີ-ທຸກ-ໆ-ທ່ານ
- On présente le résultat
 - ສະບາຍດີ (bonjour) ທຸກ (tout, tous) ໆ (?) ທ່ານ (personne)

La reconnaissance des syllabes est complexe

Exemple sur le lao

Forme générale des syllabes = **C [C] [A] V [C]**
Formes et positions problématiques des voyelles



Le nombre des syllabes peut être contraignant

■ Birman :

- Forme générale : (C ou CS) [L] V [C ou $\overset{\circ}{\text{C}}$ ou $\overset{\xi}{\text{C}}$ [$\overset{\circ}{\text{C}}$] [$\overset{\circ}{\text{C}}$ ou :]]
- Cardinaux des constituants : $|C| = 33$, $|CS| = 20$, $|L| = 15$, $|V| = 35$
- Majorant : $(33+20) \times 16 \times 35 \times (34+2) \times 2 \times 3 = \mathbf{6\ 410\ 880\ syllabes}$

■ Khmer :

- Forme générale : (C [CS [CS]] [D1] ou CS [CS]) V [C [CS] [D2]] ou VI ou L
- Cardinaux des constituants : $|C| = 33$, $|CS| = 32$, $|V| = 33$, $|VI| = 14$, $|L| = 10$
- Majorant : $(33 \times 33 \times 33 \times 4 + 32 \times 33) \times 33 \times (34 \times 33 \times 3) + 14 + 10 = \mathbf{16\ 084\ 538\ 736\ syllabes}$

■ Laotien :

- Forme générale : (C ou GC) [A] V [CF]
- Cardinaux des constituants : $|C| = 27$, $|A| = 4$, $|GC| = 36$, $|V| = 38$ et $|CF| = 8$
- Majorant : $(27+36) \times 5 \times 38 \times 9 = \mathbf{95\ 760\ syllabes}$

■ Siamois (thaï) :

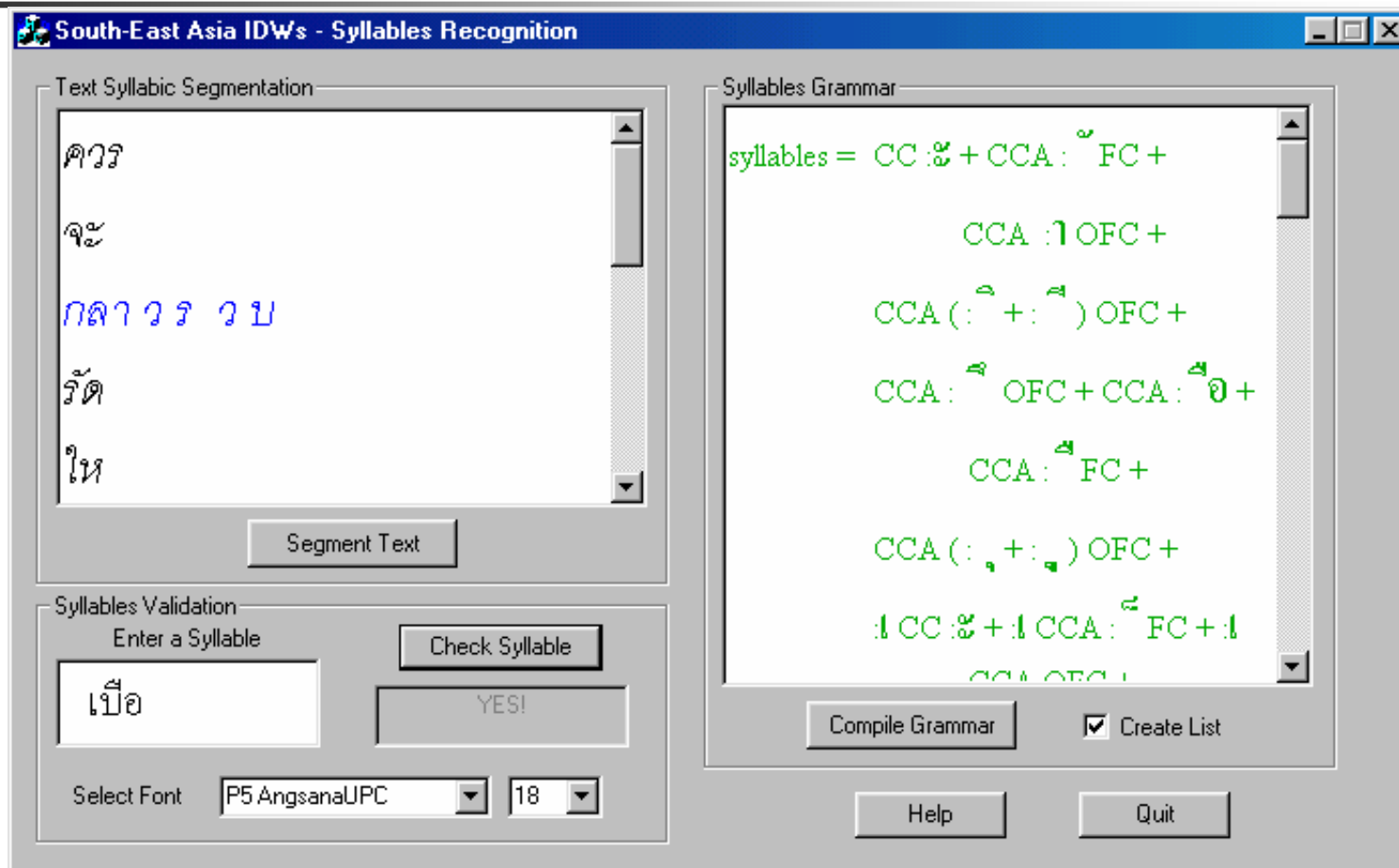
- Forme générale : (C ou GC) [A] V [CF]
- Cardinaux des constituants : $|C| = 44$, $|GC| = 140$, $|A| = 4$, $|V| = 41$, $|CF| = 38$
- Majorant : $(44+140) \times 5 \times 41 \times 39 = \mathbf{1\ 471\ 080\ syllabes}$

Syllabes : expression hors-contexte

Syllabes =	$CC : \text{z} + CCA : \text{~} CF + CCA : \text{a} CFO +$ $CCA (: \text{~} + : \text{~} + : \text{~} + : \text{~} + : \text{~} + : \text{~}) CFO +$...
CCA =	CC + CC Acc ;
CC =	GC + CI ;
GC =	:ຫ (: ງ + : ຍ + : ນ + : ມ + : ລ + : ັ + : ວ) + CI : ວ ;
CI =	: ກ + : ຂ + : ຄ + : ງ + : ຈ + : ສ + : ຊ + : ຍ + : ດ + : ຕ + : ຖ + : ທ + : ນ + : ບ + : ປ + : ຜ + : ຟ + : ພ + : ຟ + : ມ + : ຢ + : ຮ + : ລ + : ວ + : ຫ + : ອ + : ຮ + : ຫ + : ຫ ;
CFO =	CF + {} ;
CF =	: ກ + : ດ + : ບ + : ງ + : ນ + : ມ + : ຍ + : ວ ;
Acc =	: + : + : + : ;

- Appropriable par un linguiste
- Testé sur le laotien et sur le khmer

Le LSPL SYLLA



- Mise au point des automates de reconnaissance des syllabes
 - Réduit le temps de développement d'environ 80 % (60 h au lieu de 300 en moyenne)
 - Facilement appropriable
 - Utilisé pour : birman, khmer, laotien et thaï

Applications au traitement de la parole

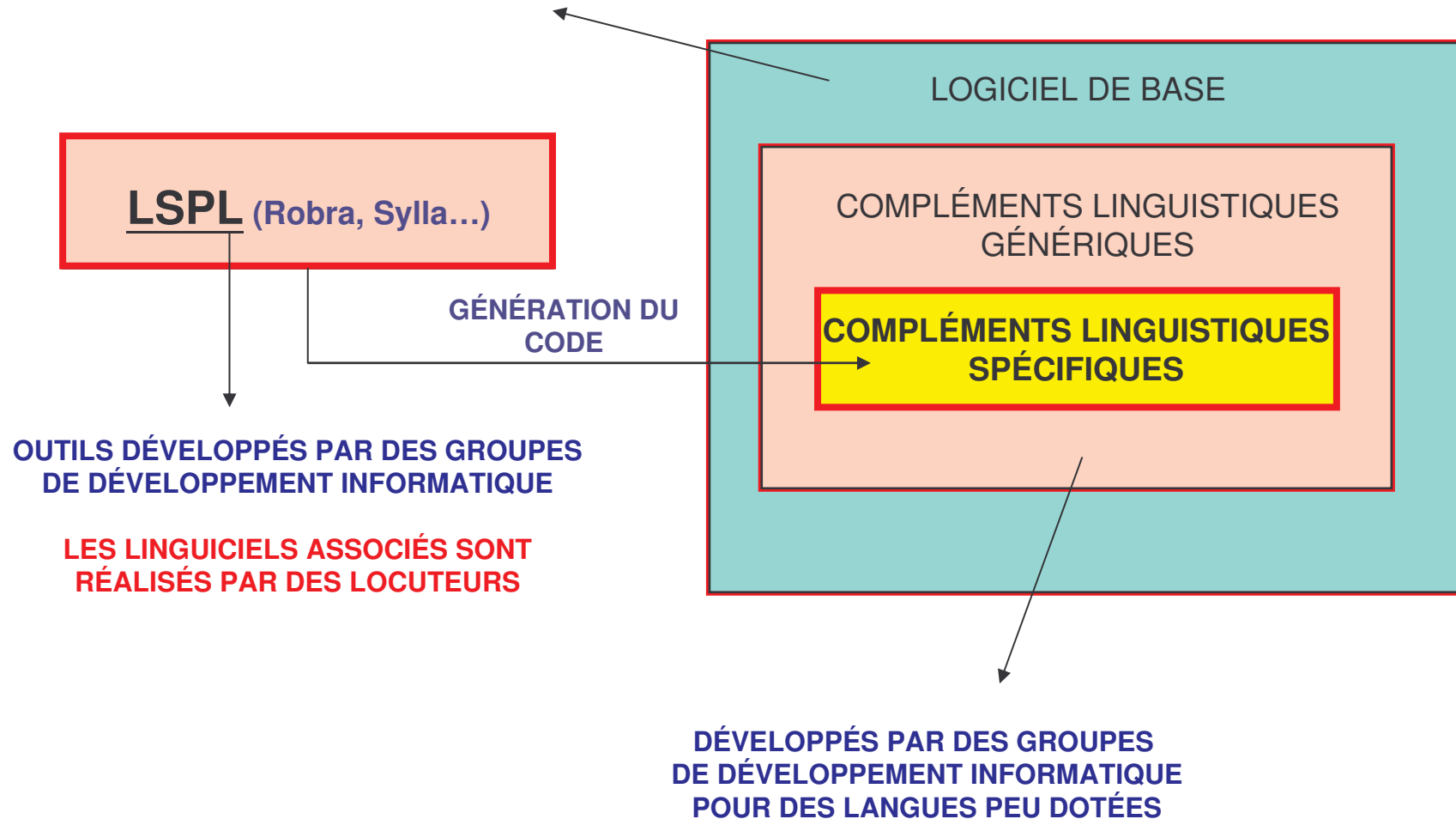
- Transcription phonétique
 - Travaux des thaïs (Thai Soundex, Thai TTS)
 - Transcriptions sur LaoSoftware.com
- Synthèse vocale et reconnaissance de la parole
 - Travaux du LIG
 - Laurent Besacier
 - Le Viet Bac
 - Sethserey Sam
 - Travaux des thaïs (Thai Soundex, Thai TTS)

Problème résiduel mais essentiel des ambiguïtés

- Exemple d'ambiguïté (thaï)
 - ผม-ตากร-ลม → ผมตากลม → ผม-ตา-กลม
 - (« je prends l'air » / « j'ai les yeux globuleux »)
- Exemples de motifs irréductibles (lao)
 - ແກງອວ(ວອວນອວ)*ໆອັງ
 - ເກວີອວ(ວ)*ບອວ(ວອວນອວ)*ໆອັງ

En conclusion : une répartition des tâches...

**DÉVELOPPÉ PAR LES GRANDS ÉDITEURS DE LOGICIELS
POUR DES LANGUES BIEN DOTÉES
DOIT PRÉSENTER DES INTERFACES STANDARDISÉES**



En conclusion : ... associée à une méthodologie

- Normaliser les interfaces entre modules
- Définir et réaliser les parties génériques
- Définir les LSPL et les outils les compilant
- Développer des linguiciels pour quelques langues, qui serviront de modèles
- Permettre la prise en charge des langues par leurs locuteurs (outils en ligne)

Site collaboratif pour un groupe de langues

The screenshot displays the GMSWare.org website interface. At the top, there is a header with a logo on the left and a multi-language navigation bar. The logo features a globe with the text 'GMSWare.org'. The navigation bar includes links for Khmer, Thai, Tai Dam, Mon, Lu, Lanna, Shan, Cham, Tham Lao, Burmese, and Lao, each followed by its respective text in that language. A search box and a 'Text Size' control are also present in the top right.

Below the header, a blue navigation bar contains 'Hide | Home' on the left and 'Edit | History | Recent Changes | Hide' on the right. The main content area is titled 'Home Page' and features a large heading: 'Welcome to GMS softWare house initiative (GMSWare.org)'. Below this heading is a section titled 'Shared vision: Preserving minority languages of Great Mekong Sub-region through cooperative framework.' followed by a paragraph of text and a link to 'Read more >>Cooperative GMS ICT Framework Case study'.

A sidebar on the left contains the following links: 'GMSWare.org', 'HomePage', 'Common', 'GMSWord', 'Languages', 'Lao', 'Dham', 'Khmer', and 'edit SideBar'. At the bottom of the sidebar, it says 'TriadSkin 3 Powered by PmWiki'.

At the bottom of the page, a Windows taskbar is visible with the system tray showing 'Chargé' and various icons.

Guider les candidats à la contribution

GMS π -Languages Computerization

I want to:

- [Develop software or participate in a collaborative development for a given language.](#)
- [Know the computerization status for a given language.](#)
- [Read an introduction to the language processing techniques.](#)

[To locate the language...](#)

=> ...select a country where the language is spoken.

Select a country

- Select a country
- Cambodia
- Laos
- Myanmar
- Thailand
- Vietnam
- Yunnan

Click on the **language**:

Language	Speakers	Family	Ethnologue
Aheu	750	Austro-Asiatic	thm
Akha	60,000	Sino-Tibetan	ahk
Bisu	1,000	Sino-Tibetan	bii
Blang	1,200	Austro-Asiatic	blr
Bru, Eastern	5,000	Austro-Asiatic	bru
Bru, Western	20,000	Austro-Asiatic	brv
Cham, Western	4,000	Austronesian	cja
Chinese, Hakka	58,800	Sino-Tibetan	hak
Chinese, Mandarin	5,880	Sino-Tibetan	cmn

Merci de votre attention

Questions ?