

# Pronunciation assessment of Japanese learners of French with GOP scores and phonetic information

Vincent Laborde<sup>1</sup>, Thomas Pellegrini<sup>1</sup>, Lionel Fontan<sup>2</sup>, Julie Mauclair<sup>1,3</sup>, Halima Sahraoui<sup>4</sup>, Jérôme Farinas<sup>1</sup>

(1) IRIT - Université de Toulouse, 31062, Toulouse, France, (2) Archean Technologies, 1899 av. d'Italie, 82000, Montauban, France

(3) Université Paris Descartes, Paris, France, (4) Octogone-Lordat - Université de Toulouse, 31058, Toulouse, France

thomas.pellegrini@irit.fr, lfontan@archean.fr, sahraoui@univ-tlse2.fr



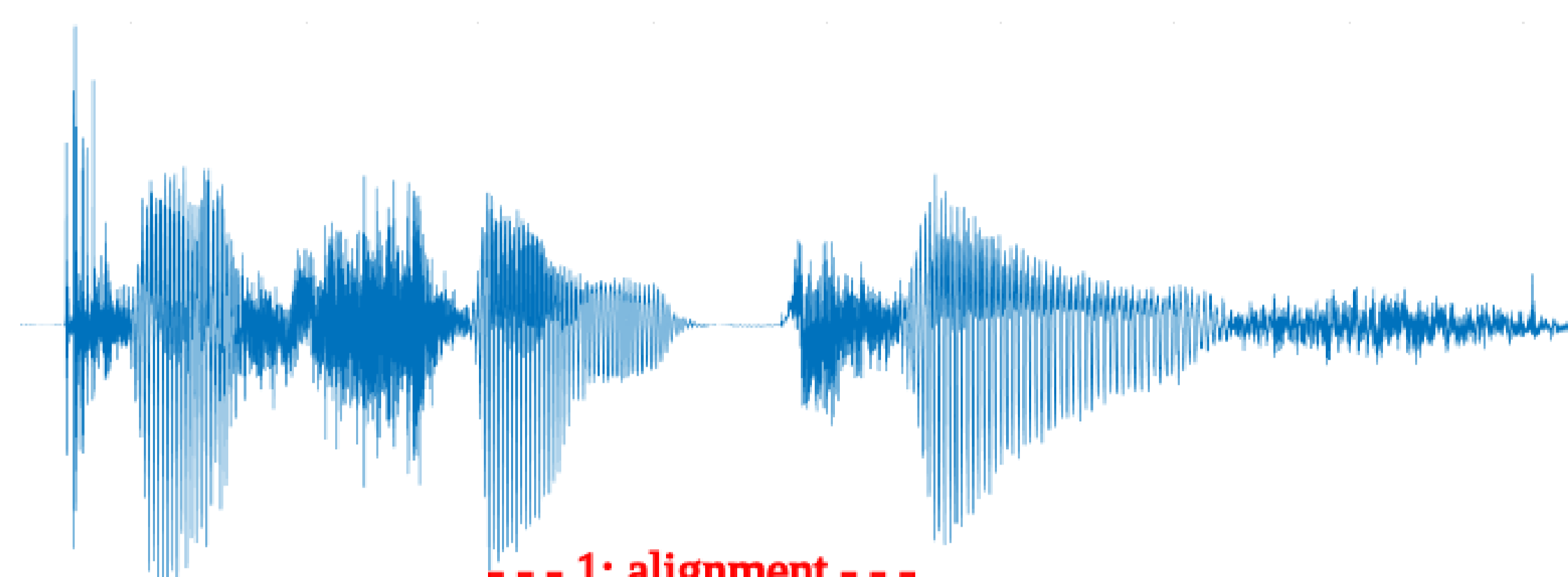
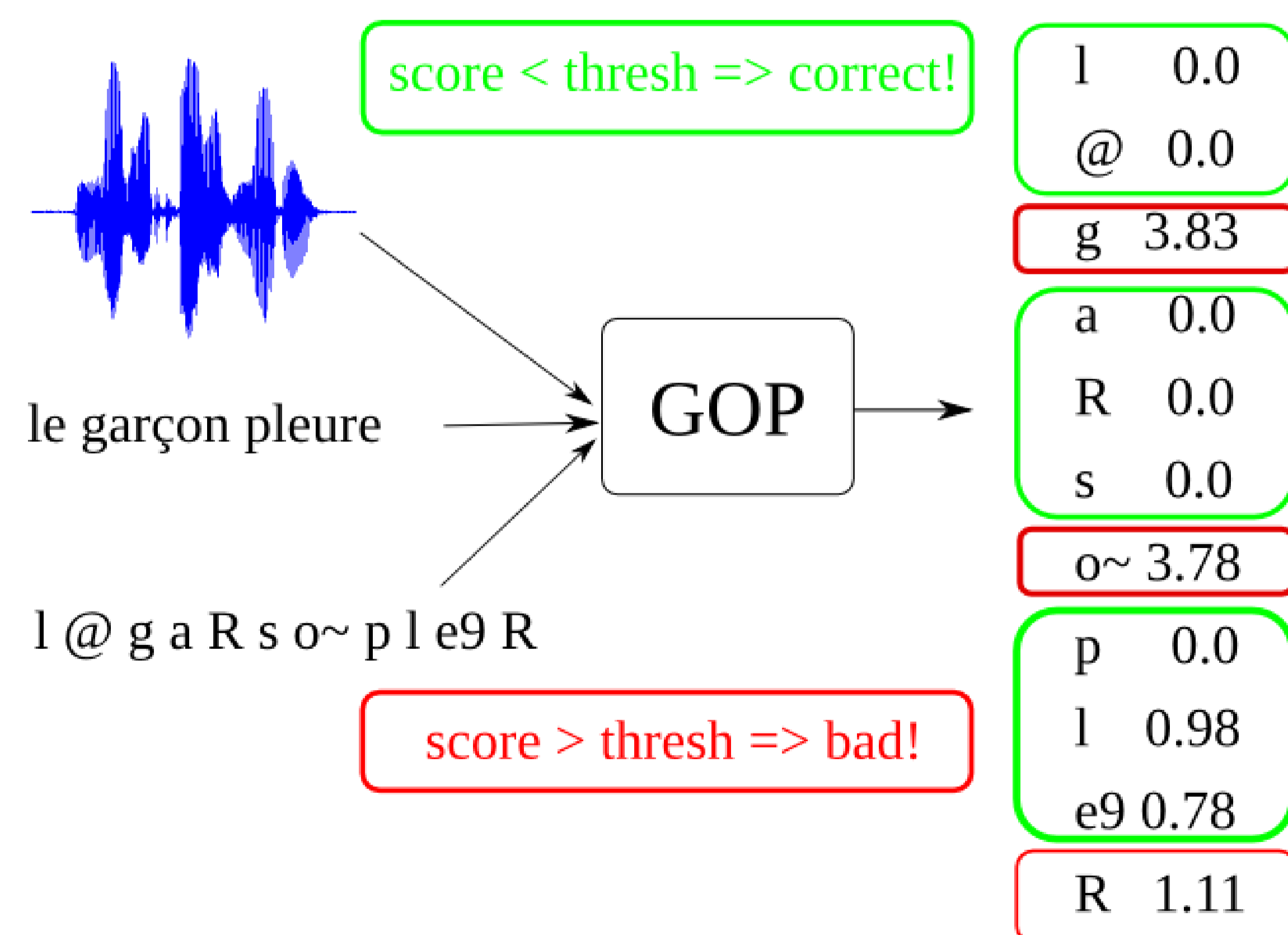
## Abstract

Automatic pronunciation assessment experiments at phone-level Japanese speakers learning French as a foreign language Three approaches compared: Goodness Of Pronunciation scores, LDA and logistic regression models Two typical errors of Japanese speakers were considered: /R/ and /v/ Best method: logistic regression when adding phonetic context information, leading to a 77.1% accuracy

## 1 Introduction

- Computer-assisted pronunciation training (CAPT) on a second language (L2)
- In this study, segmental level: one score per expected phone
- Two types of approaches:
  - Recognition scores: raw [1], likelihood ratios [3]
  - Classification methods on acoustic features: linear discriminant analysis (LDA) and alike [2]
- High accuracy is key in CAPT
- In this study:
  - Comparison between both methods + use of a classifier
  - Addition of extra phonetic and phonological features

## 2 Baseline: the forced-GOP algorithm



--- 1: alignment ---

l	g	a	R	s	o~		p	l	e9	R
		...		-702.54	-767.18	...				

--- 2: free phone recognition ---

l	k	a	R	s	a~		p	y	e2	U~
		...		-702.54	-725.51	...				

--- 3: GOP score computation ---  
 $score(o\sim) = (-725,51 + 767,18) / 11 = 3,788$

$$f\text{-GOP}(p) = \left| \log(P(O^p|p)) - \max_{j=1..J} \log(P(O^p|p_j)) \right|$$

## 3 Logistic Regression: f-GOP + LR

- Linear classifier:  $P(y = +1) = 1 / (1 + \exp(-\theta x))$
- Input: F-GOP scores and extra features
- Pros: no decision thresholds to set, feature importance given by the weights

## 4 Speech material

- BREF: read speech corpus recorded from French native speakers, 100 hours, 120 speakers
- Artificial pronunciation errors were introduced in BREF
- PHON-IM: disyllabic words repeated by FSL Japanese native speakers, 1 hour, 23 speakers

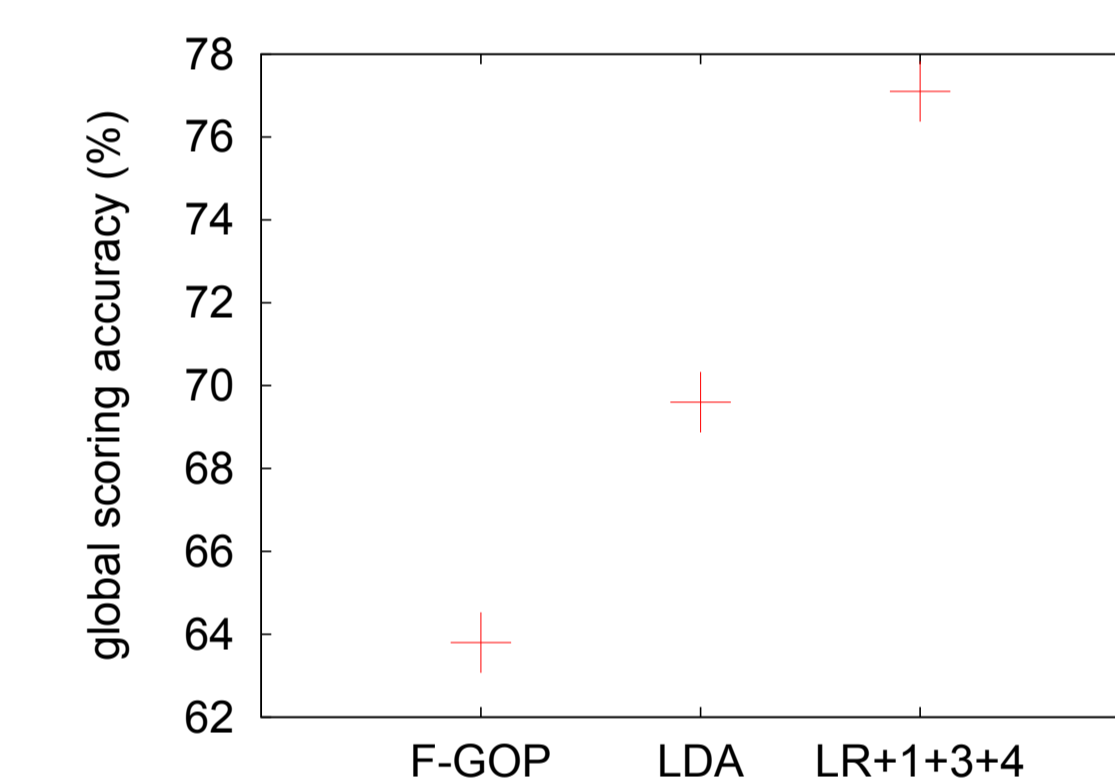
corpus	BREF		PHON-IM		Phoneme	Position		
	correct	incorrect	correct	incorrect		initial	intervocalic	final
/R/	21K	16K	215	128	/R/	47.3%	50.5%	88.3%
/v/	5K	3K	267	50	/v/	74.8%	92.6%	88.0%

## 5 Additional input features

1. Identity of the recognized phone
2. Log-likelihoods of the expected and recognized phones
3. Number of distinctive phonological features that differ between the two phones
4. Identity of the left and right phone neighbors, if any
5. Ratio between the phone duration and the duration of the middle state of the HMM

## 6 Results

- HTK, context-independent acoustic models — 39 monophones
- Scoring accuracy:  $SA = ((CA + CR) / (CA + CR + FA + FR)) \times 100$



- Global SA: f-GOP (63.8%) < LDA (69.6%) < best LR (77.1%)
- /R/ SA: LDA (62.4%) < f-GOP (68.5%) < best LR (69.1%)
- /v/ SA: f-GOP (58.7%) < LDA (77.3%) < best LR (85.8%)
- /R/ recognition: 55% as [R], 13% as [f] and 9% as a pause
- /v/ recognition: 25% as [v], 41% as [f], 1% as [b]
- Best extra features: the identity of the recognized phone, the number of distinctive phonological features, the phone neighbor identity

## 7 Conclusions

- LDA with acoustic features is better than standard f-GOP
- Best performance: f-GOP scores and phonetic / phonological features as input to LR (77% accuracy)
- Future experiments: add acoustic features to LR, deep neural networks

## References

- [1] B. Sevenster, G. de Krom, and G. Bloothoof. Evaluation and training of second-language learners' pronunciation using phoneme-based HMMs. In *Proc. STILL*, pages 91–94, Marholmen, 1998.
- [2] H. Strik, Khiet P. Truong, F. de Wet, and C. Cucchiari. Comparing classifiers for pronunciation error detection. In *Proc. INTER-SPEECH*, pages 1837–1840, 2007.
- [3] S.M. Witt. *Use of Speech Recognition in Computer-Assisted Language Learning*. Phd dissertation, University of Cambridge, Dept. of Engineering, 1999.