



IBM Research

# Large-Scale Video Semantic Filtering

**Ching-Yung Lin**

Exploratory Stream Processing Systems,  
IBM T. J. Watson Research Center

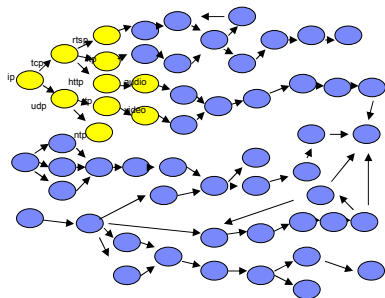
November 16, 2005

© 2004 IBM Corporation

IBM Research



## Outline -- Large-Scale Video Semantic Filtering



- Introduction and Motivation
- Prior Arts on Video Semantic Classification
- Compressed-Domain Feature Extraction
- Complexity-Accuracy Trade-Off
- Speed Up SVM Classification
- Speech-Text Complexity-Accuracy Analysis
- Conclusions
- Prospective Projects

Multimedia Semantic Routing

© 2004 IBM Corporation

## System Overview

Create **breakthrough technology** to enable

...**large-scale** production, analysis and management of

...**information and knowledge** from

...**thousands of disparate input sources**

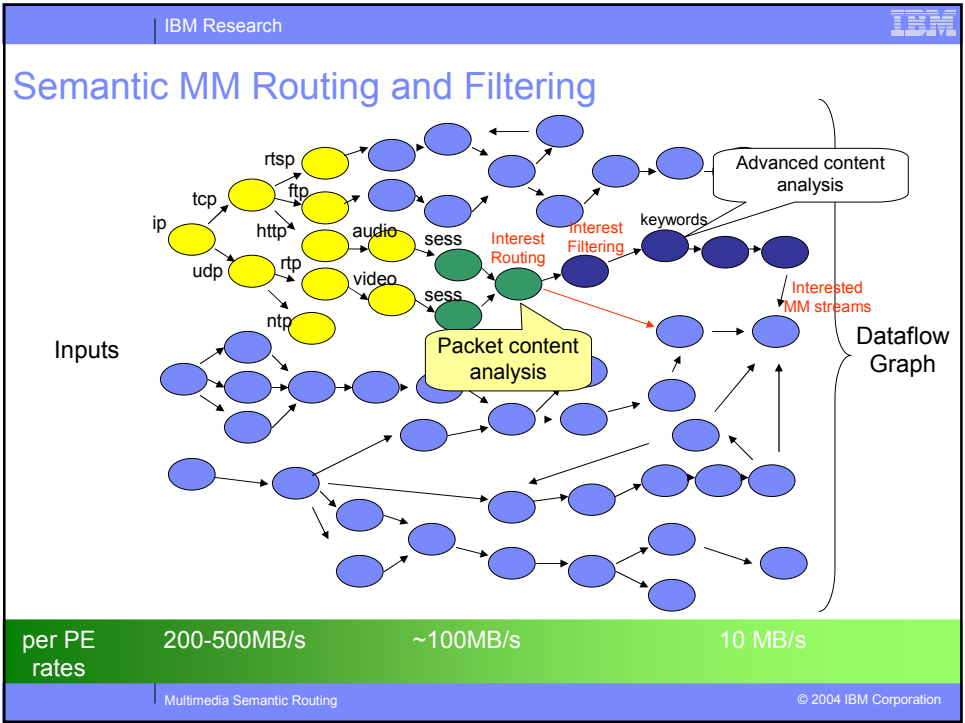
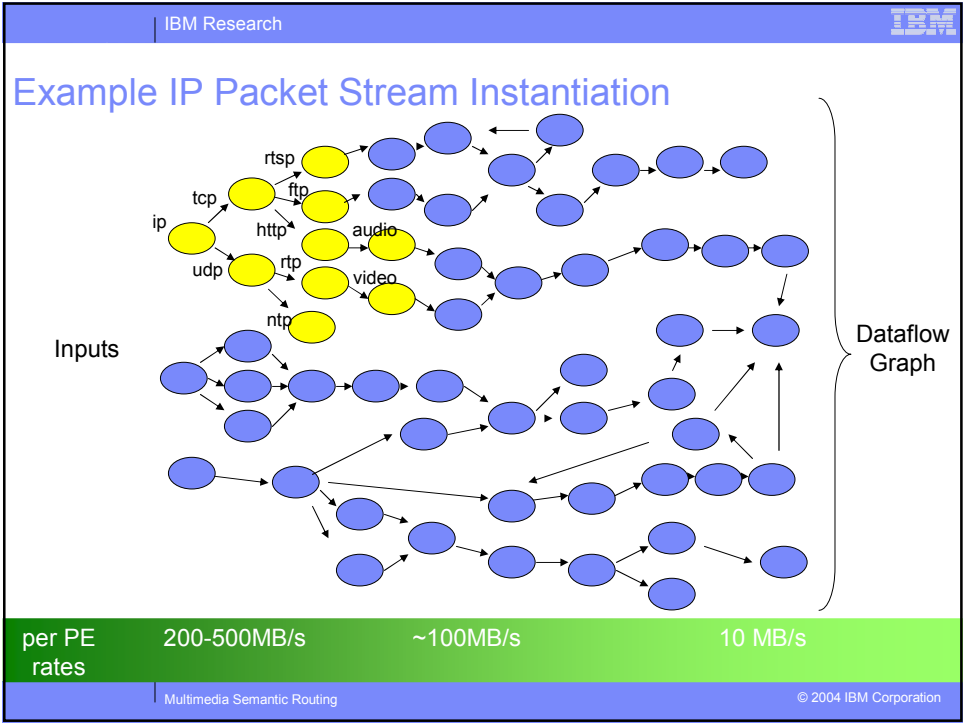
...using an **autonomic system** that

...**continuously adapts** processing to satisfy demands of

...**changing environment, content, and requests.**

## What is Large-scale?

- 10Gbit/s Continuous Feed Coming into System
- Types of Data
  - Speech, text, moving images, still images, coded application data, machine-to-machine binary communication
- System Mechanisms
  - Telephony: 9.6Gbit/sec (including VoIP)
  - Internet
    - ✓ Email: 250Mbit/sec (about 500 pieces per second)
    - ✓ Dynamic web pages: 50Mbit/sec
    - ✓ Instant Messaging: 200Kbit/sec
    - ✓ Static web pages: 100Kbit/sec
    - ✓ Transactional data: TBD
  - TV: 40Mb/sec (equivalent to about 10 stations)
  - Radio: 2Mb/sec (equivalent to about 20 stations)



IBM Research

## Objective

- Input data X – Queries q – Resource R**
  - $Y(X | q)$ : Relevant information
  - $Y'(X | q, R) \in Y(X | q)$ : Achievable subset given R
- Configurable Parameters of Processing Elements to maximize relevant information:**
  - $Y''(X | q, R) > Y'(X | q, R)$ ,  
with resource constraint.
- Required resource-efficient algorithms for:**
  - Classification, routing and filtering of signal-oriented data: (audio, video and, possibly, sensor data)

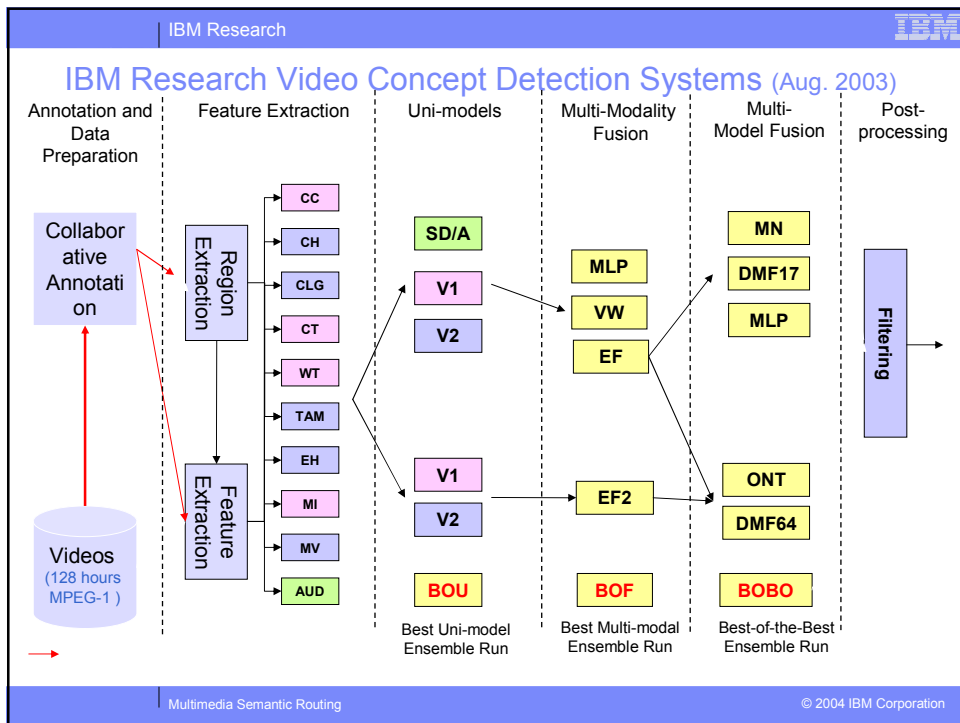
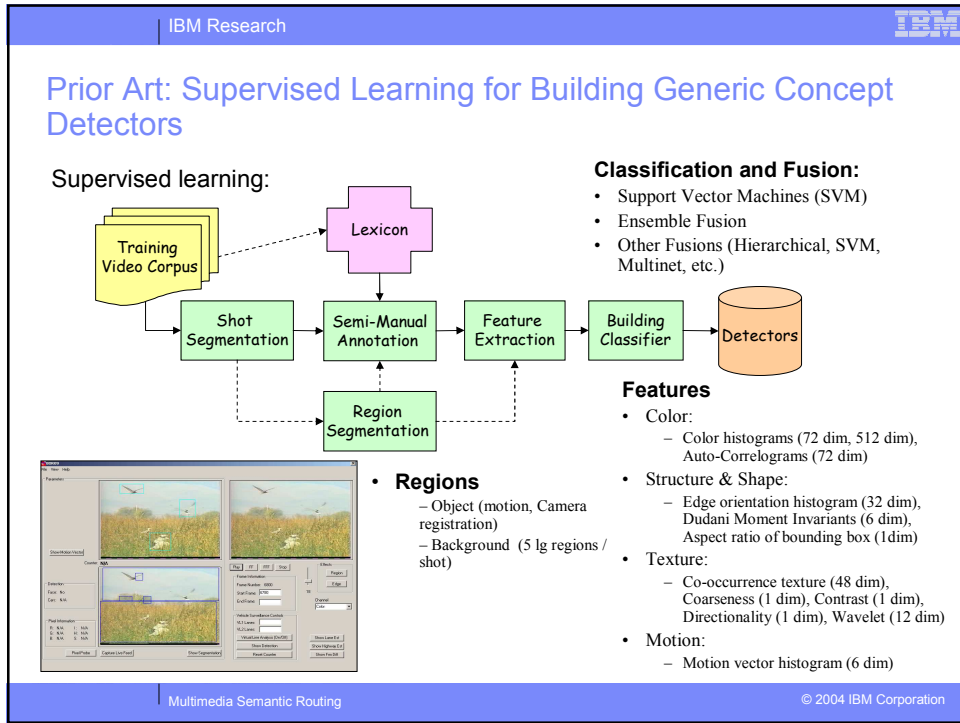
Multimedia Semantic Routing © 2004 IBM Corporation

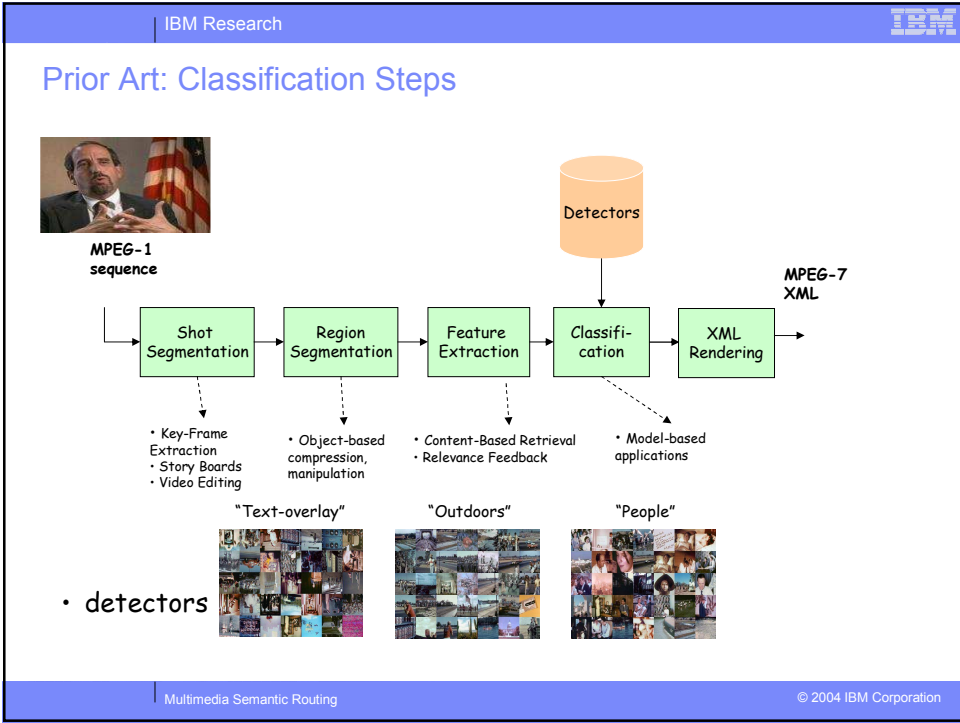
IBM Research

## Video Concept Classification/Routing

- Objective:**
  - Build Moderate Amount of Real-Time Concept Classifiers, e.g., Human, Outdoors, Face, Indoors, etc.
- Test-bed Corpus:**
  - NIST TRECVID 2003 corpus: 128 hours of MPEG-1 videos, 62 hours of them were manually annotated by 23 worldwide groups using IBM *VideoAnnEx* Collaborative Annotation System.
  - 38 hours of video for training; 24 hours (partitioned to 6-, 6-, 12-hour sets) for internal evaluation; 65 hours for NIST benchmarking.

Multimedia Semantic Routing © 2004 IBM Corporation





IBM Research

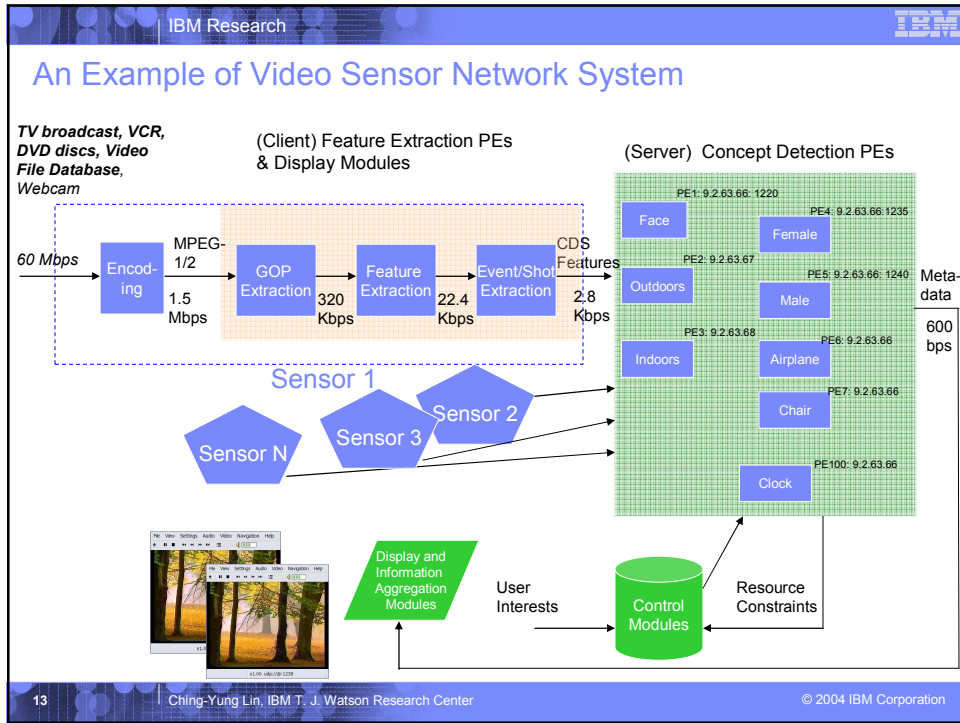
## Novel Visual Semantic Concept Filters

- 100 visual filters were built (v.s. 64 in 2003).
- Concepts are arranged hierarchically.
- Based on novel compressed-domain sliced features.
- Neither segmentation nor fusion is used.

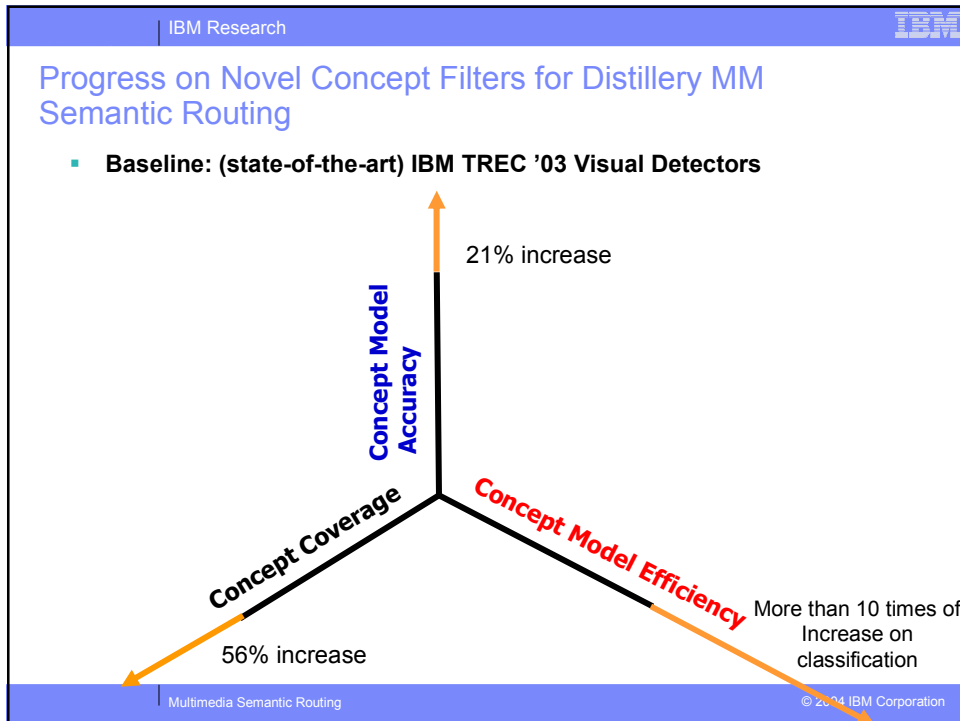
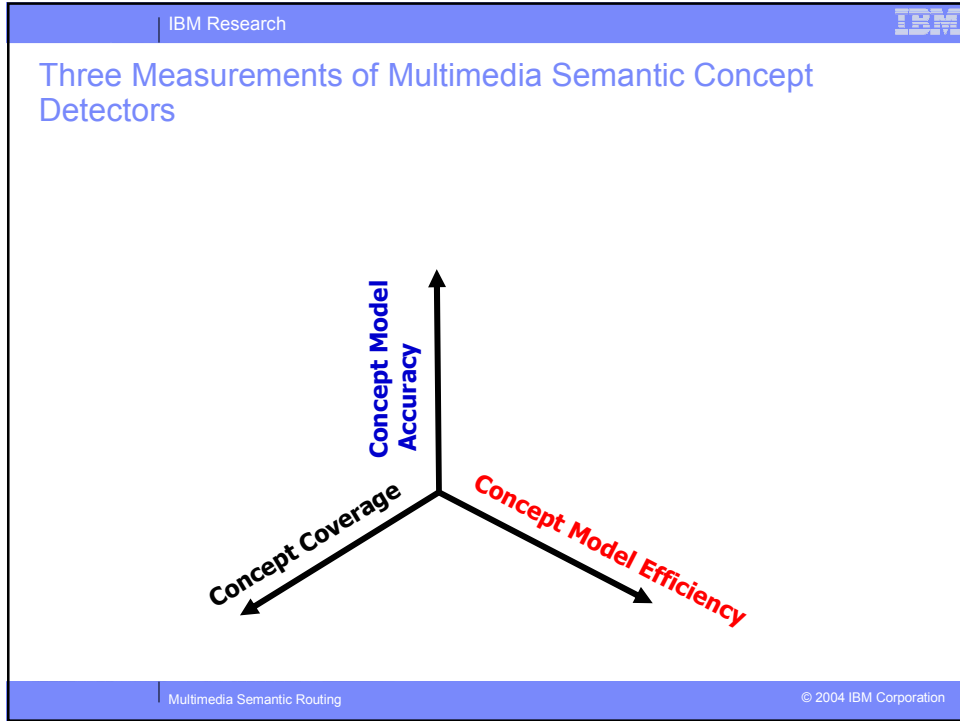
The screenshot shows a hierarchical interface for visual semantic concept filters. It is divided into three main sections: Key Frame, Event, and Scene. The Key Frame section shows a red car. The Event section lists various actions like Person\_Action, Monologue, News\_Subject\_Mor, etc. The Scene section is further divided into Indoors, Non-Studio\_Setting, Studio\_Setting, etc. The Object section lists various objects like Animal, Chicken, Cow, etc.

Multimedia Semantic Routing

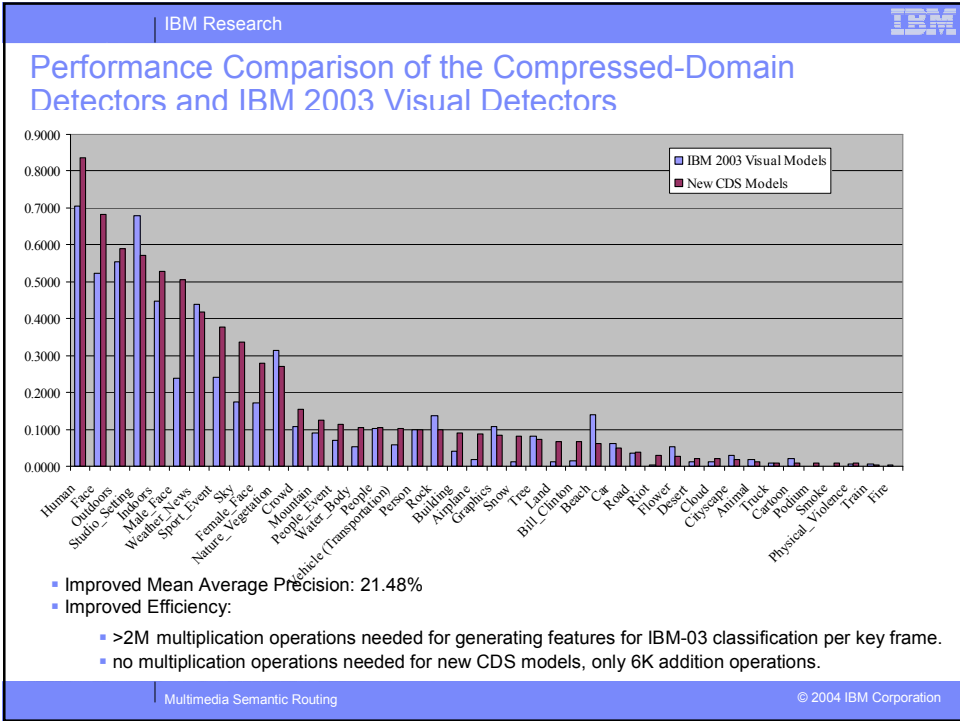
© 2004 IBM Corporation



- IBM Research
- ## Main Process
- ❑ This feature set is extracted as follows:
    1. Parse the MPEG-1/2 packets and get the beginning of an I-frame or the closest I-frame of a pre-specified shot keyframe.
    2. Using the VLC maps to map variable-length codes to the DCT-domain coefficients.
    3. Within an MPEG slice or a union of slices, truncate selected DCT coefficients and calculate the histogram of these DCT coefficients.
    4. Form a feature vector of the frame based on the histogram coefficients with multiple slices.
  - ❑ In the above procedure, we can see that no multiplication operation is required to get these feature vectors. Only addition is needed for getting the histogram.
  - ❑ In a typical situation, we partition a frame into three slices, and use 3 DCT coefficient histogram (1 DC and 2 lowest frequency AC coefficients) on all YCbCr domains. This forms a 576-dimensional feature vector.
  - ❑ After feature extraction, SVM is used to train models and classification.
- 14 Ching-Yung Lin, IBM T. J. Watson Research Center © 2004 IBM Corporation







IBM Research

## Demo -- Novel Semantic Concept Filters

- <http://www.research.ibm.com/VideoDIG>
- E.g.:

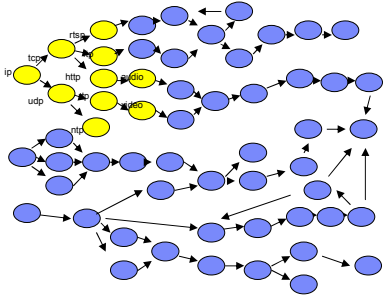
Multimedia Semantic Routing © 2004 IBM Corporation

IBM Research

INTERNATIONAL BUSINESS MACHINES CORPORATION  
 500 WEST MONTELEONE AVENUE  
 CHICAGO, ILLINOIS 60611-1000  
 WWW.IBM.COM

## Outline -- Large-Scale Video Semantic Filtering

- Introduction and Motivation
- Prior Arts on Video Semantic Classification
- Compressed-Domain Feature Extraction
- Complexity-Accuracy Trade-Off
- Speed Up SVM Classification
- Speech-Text Complexity-Accuracy Analysis
- Conclusions
- Prospective Projects



Multimedia Semantic Routing

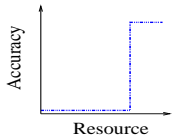
© 2004 IBM Corporation

IBM Research

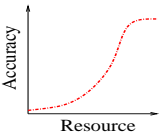
INTERNATIONAL BUSINESS MACHINES CORPORATION  
 500 WEST MONTELEONE AVENUE  
 CHICAGO, ILLINOIS 60611-1000  
 WWW.IBM.COM

## Complexity Reduction Introduction

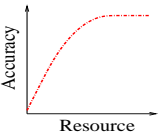
- **Objective: Real-time classification of instances using Support Vector Machines (SVMs)**
- **Computationally efficient and reasonably accurate solutions**
- **Techniques capable of adjusting tradeoff between accuracy and speed based on available computational resources**



SVM



Objective



Achieved

Multimedia Semantic Routing

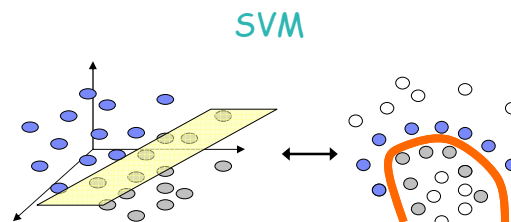
© 2004 IBM Corporation

## SVMs

- **Developed over the past decade**
- **Strong theoretical background**
  - Attempt to achieve Structural Risk Minimization
  - Minimize upper bound on generalization error instead of MSE
- **Impressive empirical results, handles noisy data**
- **Effective over diverse domains**
  - Text (character/word based), Image, Surveillance (Video/Network traffic), Bio (Gene/Protein/Disease) etc.

## SVM formulation

- **Given :**
  - Training instances  $\{x_i\}$  with labels  $y_i$
- **Objective :**
  - Find maximum margin hyperplane separating positive and negative training instances



## Salient Features

- **Formulation:**

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

$$s.t. \quad y_i(\mathbf{w} \cdot \phi(x_i) + b) \geq 1$$

- **Optimization problem with quadratic objective function, linear constraints (QPP)**
- **Hyperplane separator in *projected* space**
- **Implicit projection of instances**

## Kernel Projection

- **Implicit projection to feature space**

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$

- **Measures similarity between instances in projected space**
- **Examples:**
  - Gaussian :  $\text{Exp}(-\gamma \|x_i - x_j\|^2)$
  - Laplacian :  $\text{Exp}(-\gamma \|x_i - x_j\|)$
  - Polynomial :  $(1 + x_i \cdot x_j)^p$

## Decision

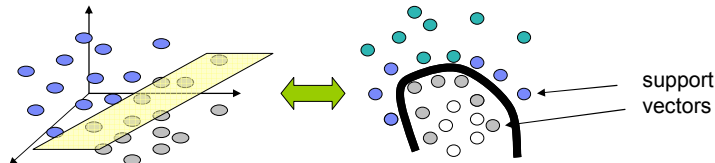
- **Score of unseen instance**  $u_j : w \cdot \phi(u_j)$

- **In terms of Lagrangian multipliers**

$$\sum_i \alpha_i y_i k(x_i, u_j)$$

- **Computational Cost :  $O(n_{sv}d)$** 
  - $n_{sv}$ : Number of support vectors
  - $d$  : Dimensionality of each data instance

## Complexity Analysis on SVM classifiers



- **Support Vector Machine**
  - Largest margin hyperplane in the projected feature space
  - With good kernel choices, all operations can be done in low-dimensional input feature space

- SVM Classifier:

$$f(x) = \sum_{i=1}^S a_i \cdot k(x, x_i) + b$$

where  $S$  is the number of support vectors,  $k(\cdot, \cdot)$  is a kernel function. E.g.,  $k(x, x_i) = e^{-\gamma \|x - x_i\|^2}$

- Complexity  $c$ : operation (multiplication, addition) required for classification

$$c \propto S \cdot D$$

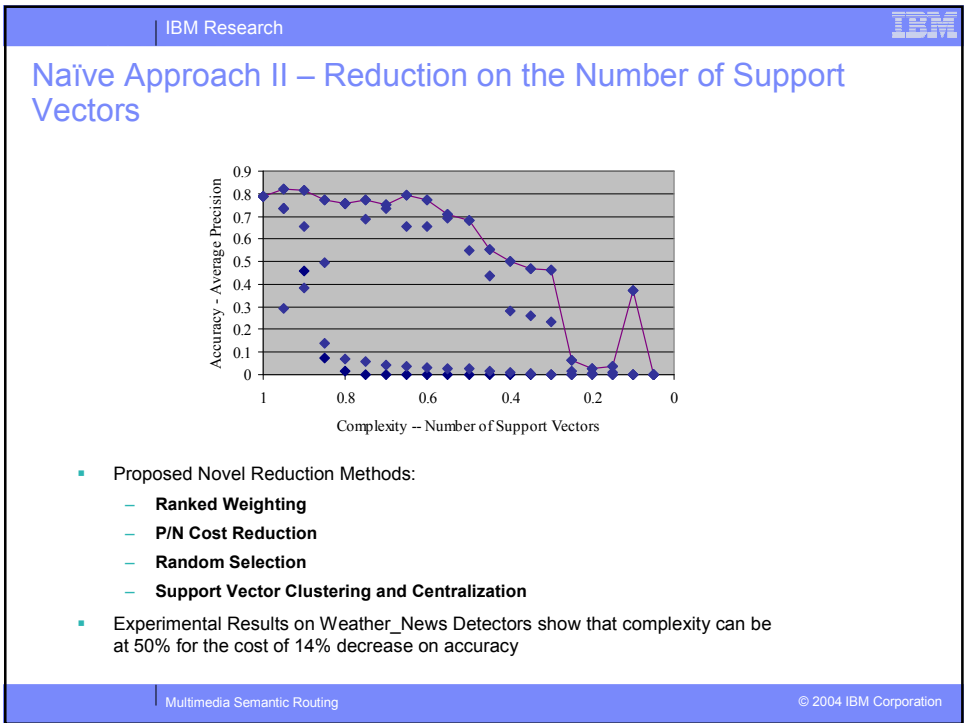
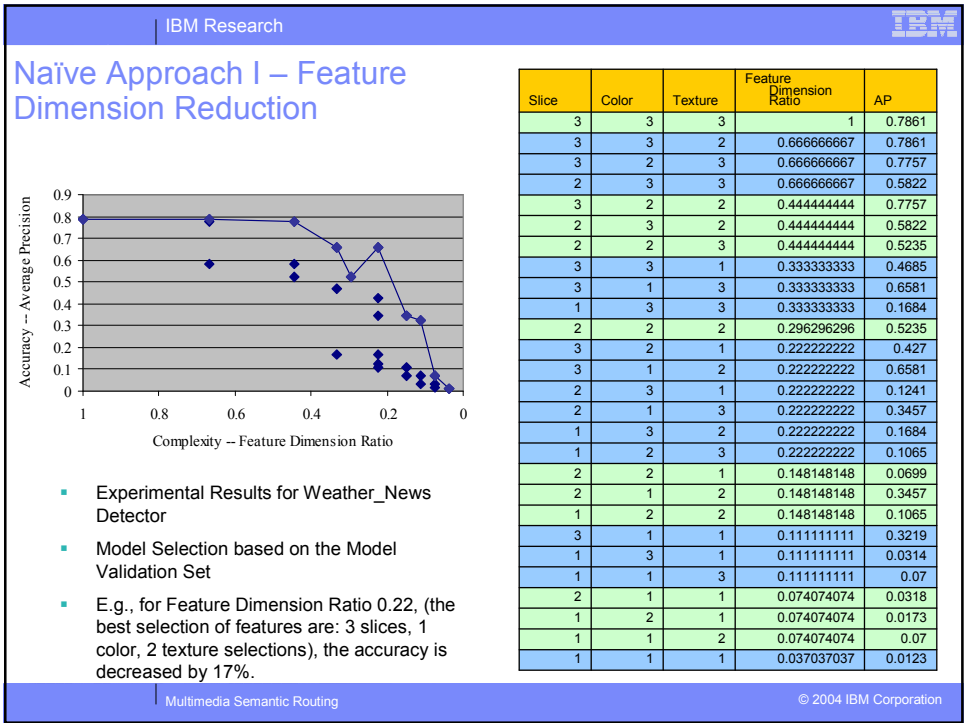
where  $D$  is the dimensionality of the feature vector

## Problems

- **Number of support vectors grows quasi-linearly with size of training set [Tipping 2000]**
- **Inner product with each support vector of dimensionality  $d$  expensive**
  - Example TREC2003
    - Human : 19745 support vectors
    - Face : 18090
- **High data rates(10Gbits/sec) means large number of abandoned data**

## Example

- **Processing Power 1 Ghz**
- **10000 support vectors**
- **1000 / 2 features per instance**
- **Order of at least  $10^7$  operations required per stream per sec**
- **Translates to less than 100 instances evaluated per sec with only one classifier**



## Approach Overview

- **Reduced set methods**
- **Focus on support vectors and weights already obtained (Same SVM formulation)**
- **Accommodate effect due to all support vectors in approximation**
- **Amenable to**
  - multiple kernels
  - $L_1, L_2$  norms
- **Approximate support vector scoring with consideration for overall error**

## Approaches

- **Weighted Clustering Approach:**
  - Approximate effect of support vector(s) by (possibly hypothetical) instance with weight adjustment
- **Error Radius Approach:**
  - Determine *error radius* and approximate effect outside same



## Weighted Clustering Approach

- **Basic steps**

- Cluster support vectors
- Use cluster center as representative for all support vectors in cluster
- Determine scalar weight associated with each cluster center
- Use only cluster centers to score new instances

## Weighted Clustering Approach (contd.)

- **Clustering**

- Cluster positive and negative support vectors separately
- Ratio of +ve clusters to -ve clusters same as ratio of +ve support vectors to -ve support vectors
- Attempts to avoid *bloated* clusters because of imbalance
- Manual decision on number of clusters

## Weighted Clustering Approach (contd.)

- **Cluster center weight**

- Given :
  - Support vector  $\mathbf{x}_i$  in cluster with cluster center  $\mathbf{x}_c$
  - Weight  $\alpha_i$  associated with  $\mathbf{x}_i$
- $\Delta_i$  Squared distance between  $\mathbf{x}_i$  and  $\mathbf{x}_c$
- Unknown instance at squared distance
  - $d$  from  $\mathbf{x}_c$
  - $d - \delta_i$  from  $\mathbf{x}_i$

## Cluster center weight (contd.)

- **Difference in scores**

$$\gamma_i \exp\left(-\frac{d}{p}\right) - \alpha_i \exp\left(-\frac{(d-\delta_i)}{p}\right)$$

- **Optimal**  $\gamma_i = \alpha_i \exp\left(\frac{\delta_i}{p}\right)$

- **Problem:**

- Determined  $\gamma_i$  specific to instance
- Impossible to obtain speedup if  $\gamma_i$  computed on a per-instance basis

- **Solution:**

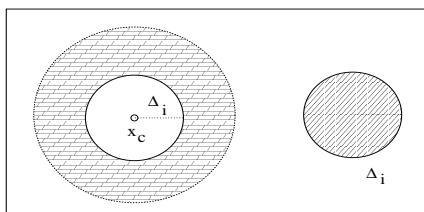
- Find  $\gamma_i$  minimizing error on average

### Cluster center weight (contd.)

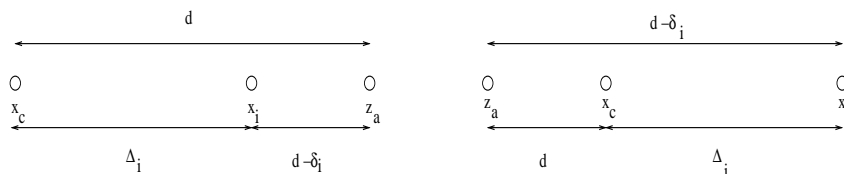
- Choose  $\gamma_i$  minimizing square of difference in scores over all  $\delta_i$  and  $d$
- Sub-cases :

$$d \geq \Delta_i$$

$$d < \Delta_i$$



### Cluster center weight (contd.)

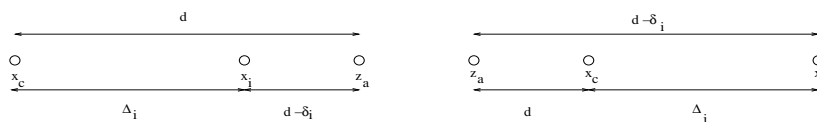


- Case 1: Distance to unseen instance  $d \geq \Delta_i$

$$\int_{\Delta_i}^{\infty} \int_{-\Delta_i}^{\Delta_i} (\gamma_i \exp(-\frac{d}{p}) - \alpha_i \exp(-\frac{(d-\delta_i)}{p}))^2 d\delta_i dd.$$

- Optimal 
$$\gamma_{i_{upper}} = \alpha_i p \frac{\exp(\frac{\Delta_i}{p}) - \exp(-\frac{\Delta_i}{p})}{2\Delta_i}$$
- Minima verification:  $2\Delta_i p \exp(-\frac{2\Delta_i}{p}) \geq 0$
- Optimal  $\gamma_{i_{upper}}$  obtained for  $d \geq \Delta_i$

## Cluster center weight (contd.)



- **Case 2:**  $d < \Delta_i$

$$\int_0^{\Delta_i} \int_{\Delta_i - 2d}^{\Delta_i} (\gamma_i \exp(-\frac{d}{p}) - \alpha_i \exp(-\frac{(d+\delta_i)}{p}))^2 \mathbf{d}\delta_i \mathbf{d}d$$

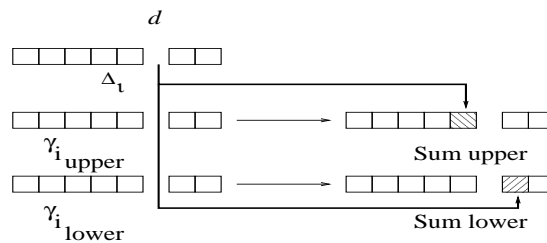
- **Optimal**  $\gamma_{i_{lower}} = \frac{\alpha_i \exp(-\frac{\Delta_i}{p}) [\frac{2\Delta_i}{p} - 1 + \exp(-\frac{2\Delta_i}{p})]}{1 - \exp(-\frac{2\Delta_i}{p}) - \frac{2\Delta_i}{p} \exp(-\frac{2\Delta_i}{p})}$
- **Minima:**  $\exp(\frac{2\Delta_i}{p}) > (1 + \frac{2\Delta_i}{p})$
- **Optimal**  $\gamma_{i_{lower}}$  obtained for  $d < \Delta_i$

## Using the weights

- **For every support vector in cluster**
  - Distance  $\Delta_i$  known
  - Two weights computed
- **Cumulative effect of all support vectors in clusters additive**
  - $\Delta_i$  because of various support vectors added up at center to simulate effect of all support vectors
- $\Delta_i$  sorted, weight arrays rearranged

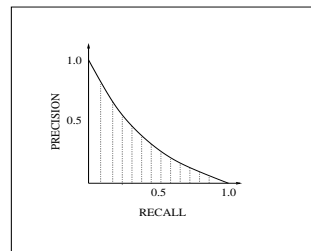
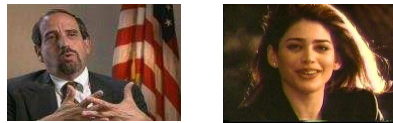
## Using the weights (contd.)

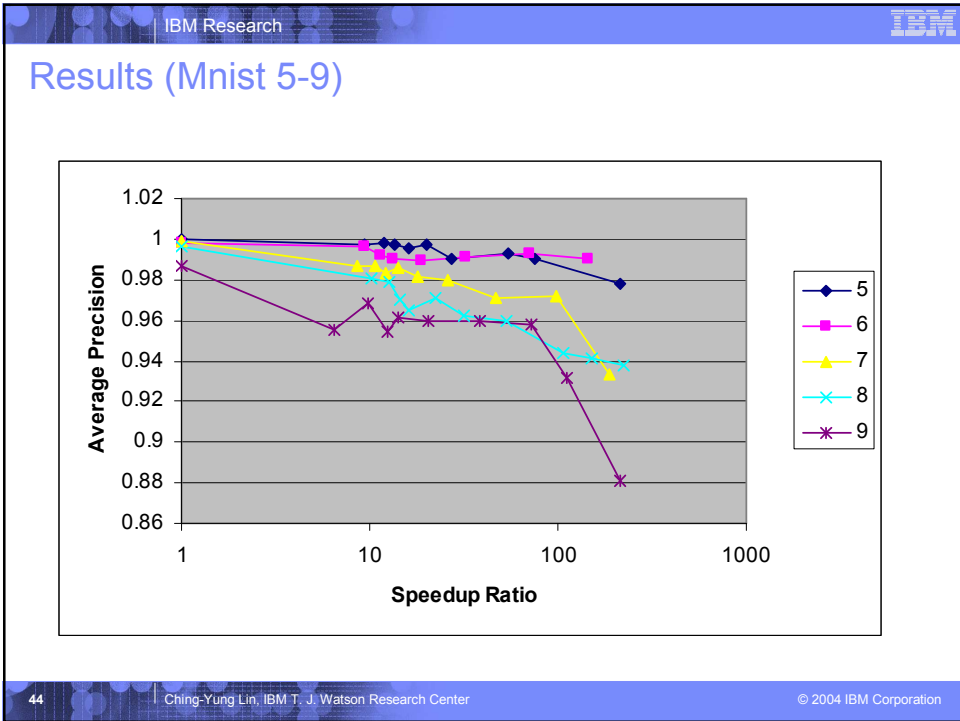
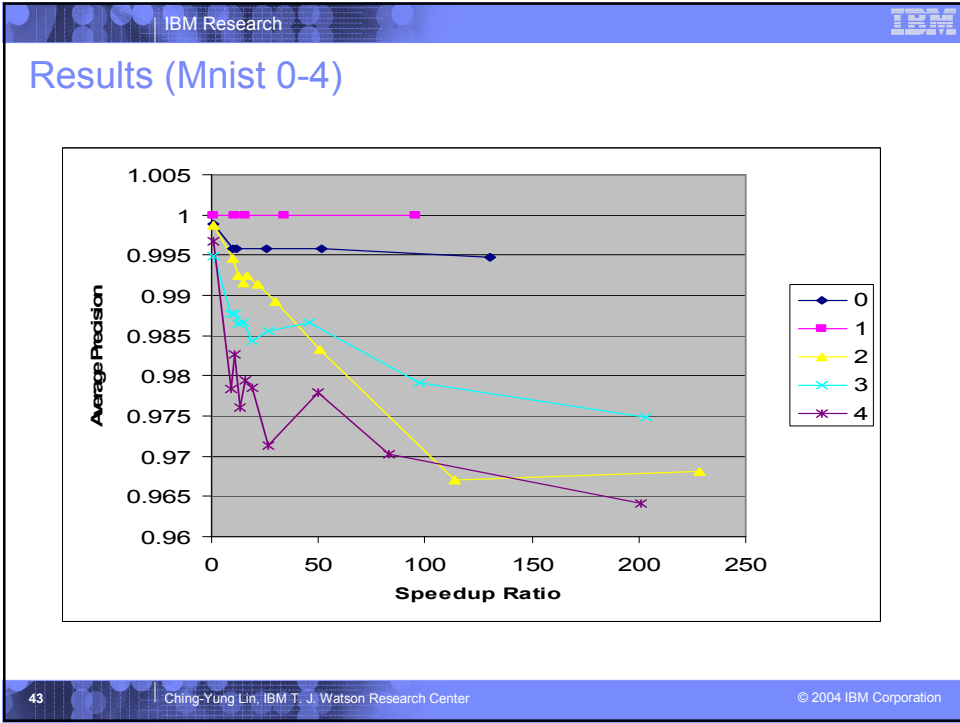
- **Given distance of unseen instance  $d$** 
  - Binary search on  $\Delta_i$  array
  - For support vectors with  $d$  less than  $\Delta_i$ ,  $\gamma_{i_{upper}}$  used,  $\gamma_{i_{lower}}$  for others
- **Sum associated with cluster center not computed per-instance**
- **Pre-computed at every position of sorted array**

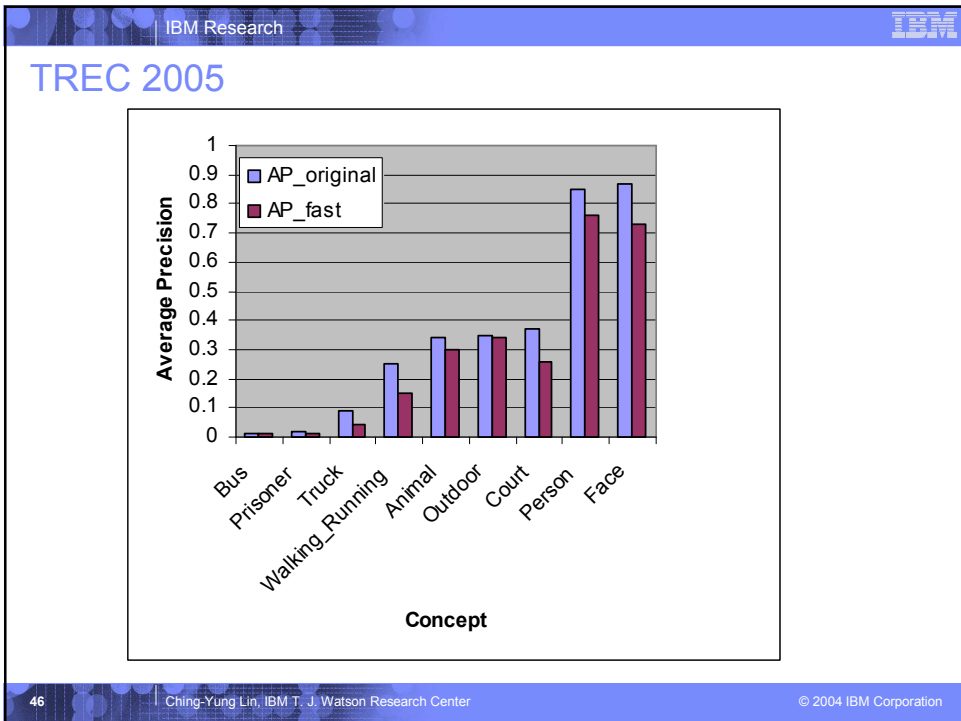
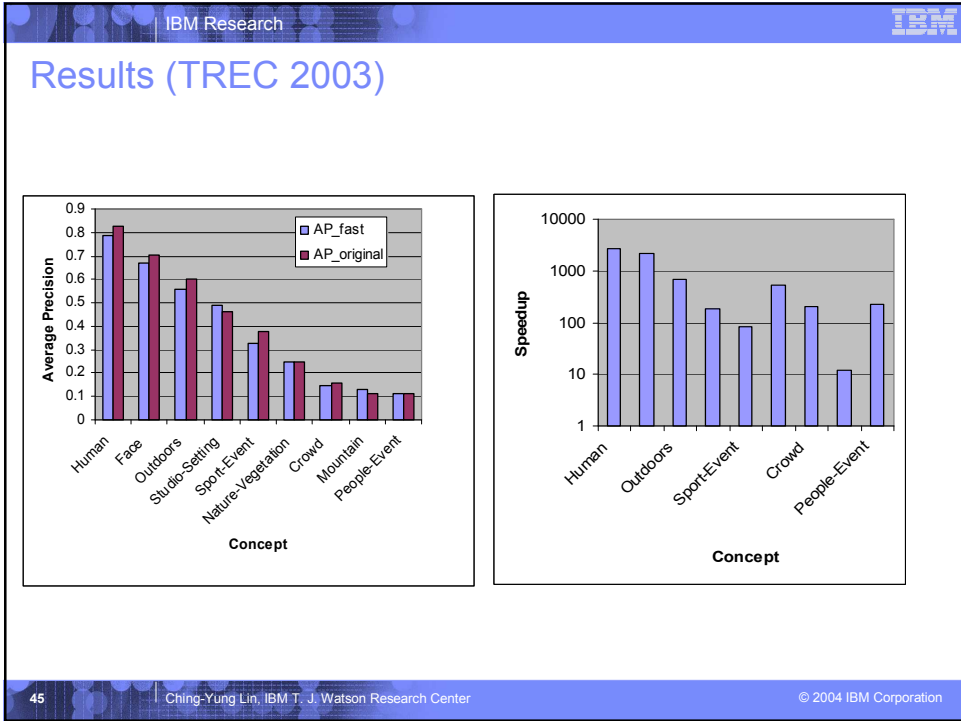


## Experiments

- **Datasets**
  - TREC video datasets (2003 and 2005)
    - 576 features per instance
    - > 20000 test instances overall
  - MNist handwritten digit dataset (RBF kernel)
    - 576 features
    - 60000 training instances, 10000 test instances
- **Performance metrics**
  - Speedup achieved over evaluation with all support vectors
  - Average precision achieved







IBM Research

## Summary of Complexity Reduction

- ❑ Techniques presented demonstrate reasonable performance in terms of both speedup and average precision over multiple concepts in datasets
- ❑ Speedups
  - MNist : All concepts at least 50 times faster with AP within 0.04 of original
  - TREC 2003: Eight out of nine concepts speedup greater than 80 times with AP within 0.05 of original
  - TREC 2005: APs in some cases along with speedup respectable
- ❑ APs of most concepts close to original APs

47 Ching-Yung Lin, IBM T. J. Watson Research Center © 2004 IBM Corporation

IBM Research

## VideoDIG Demo

– configurable and scalable video semantic concept detection/filtering

(Client) Feature Extraction PEs & Display Modules

(Server) Concept Detection PEs

TV broadcast, VCR, DVD discs, Video File Database, Webcam

60 Mbps

Encoding

MPEG-1/2

1.5 Mbps

320 Kbps

22.4 Kbps

2.8 Kbps

600 bps

Meta-data

Face

Female

Outdoors

Male

Indoors

Airplane

Chair

Clock

PE1: 9.2.63.66: 1220

PE2: 9.2.63.67

PE3: 9.2.63.68

PE4: 9.2.63.66: 1235

PE5: 9.2.63.66: 1240

PE6: 9.2.63.66

PE7: 9.2.63.66

PE100: 9.2.63.66

User Interests

Control Modules

Resource Constraints

On/Off

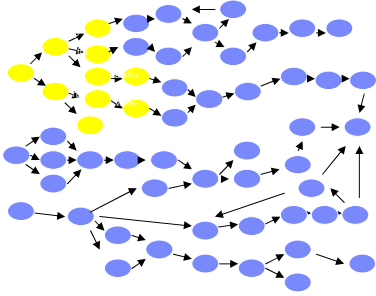
48 Ching-Yung Lin, IBM T. J. Watson Research Center © 2004 IBM Corporation



IBM Research

## Outline – Large-Scale Video Semantic Filtering

- Introduction and Motivation
- Prior Arts on Video Semantic Classification
- Compressed-Domain Feature Extraction
- Complexity-Accuracy Trade-Off
- Speed Up SVM Classification
- **Speech-Text Complexity-Accuracy Analysis**
- **Conclusion**
- **Prospective Projects**



49 Ching-Yung Lin, IBM T. J. Watson Research Center © 2004 IBM Corporation

IBM Research


## Other Issues: Speech and Text-based Topic Detection

- **Unsupervised learning from WordNet:**

Query Concept  
Weather News

WordNet

weather, weather condition, atmospheric condition  
=> cold weather  
=> fair weather, sunshine  
=> hot weather



Weather condition  
Atmospheric condition  
Cold weather  
Fair weather

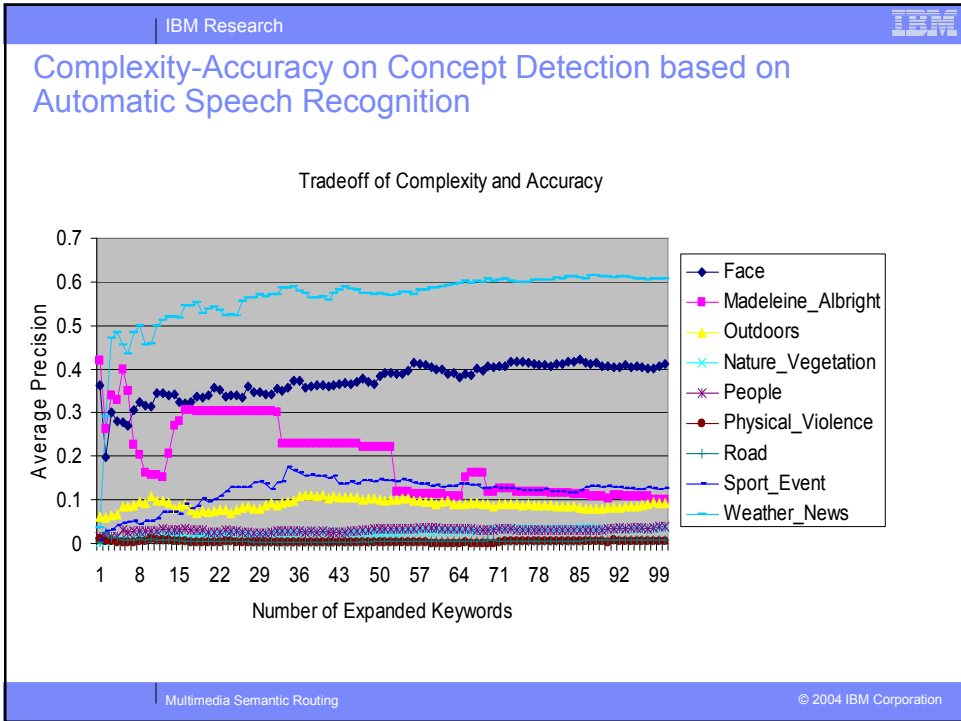
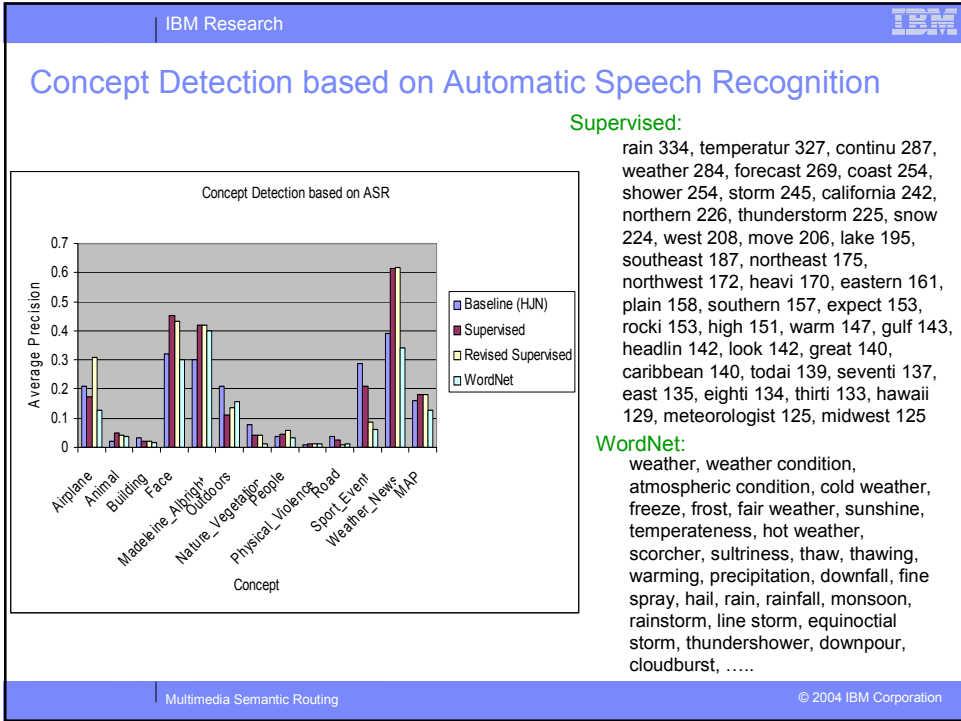
Weather forecast  
Warm weather

Speech-based Retrieval

Airplane	Animal	Building	Weather news
airplane	animal	building	Weather
aeroplane	beast	edifice	atmospheric
plane	brute	walk-up	cold weather
airline	critter	butchery	freeze
airbus	darter	apart build	frost
twin-aisle	peeper	tenement	sunshine
airplane	creature	architecture	temper
amphibian	Fauna	call center	scorcher
biplane	microorganism	sanctuary	sultry
passenger	poikilotherm	Bathhouse	thaw
fighter			warm
aircraft			downfall
			hail
			rain

Example of Topic - Related Keywords

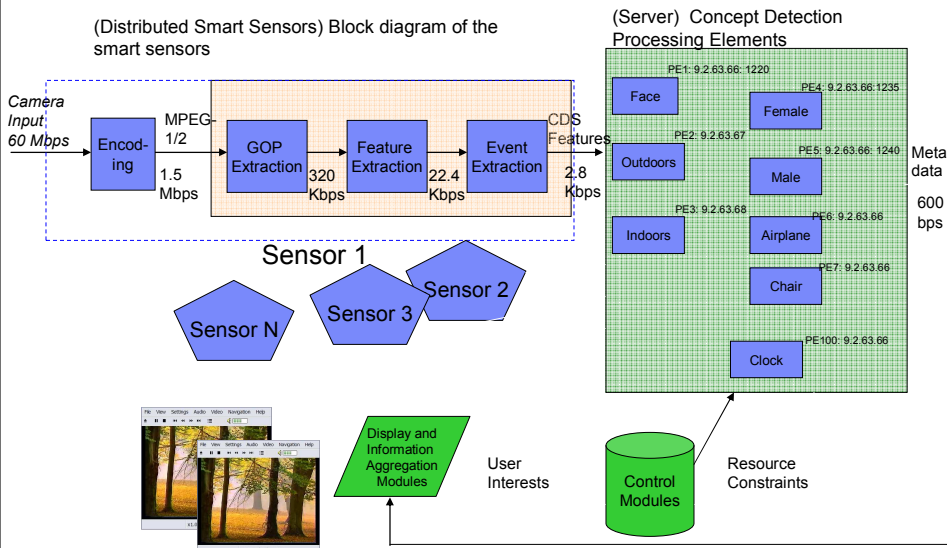
Multimedia Semantic Routing © 2004 IBM Corporation



## Conclusion

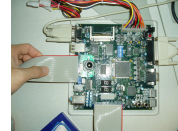
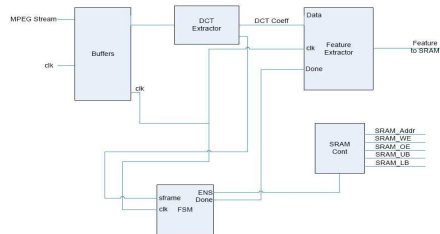
- **Novel Challenge:**
- **Current Status:**
  - Classifying 40 video streams of 100 concepts in real-time by using a Pentium III PC.
  - SVM classification can be speed up in a factor of 10s to 1000s with average precision of 90% or above of the original.
- **Future Works:**
  - Learning and Classifying in noisy environment
  - Frameworks for other types of media (e.g., speech, email, web access, ....) and other types of machine learning/pattern recognition framework (e.g., dynamic Bayesian Net..)

## Revisited Video Semantic Filtering Framework



## Projects Available (Spring 2006)

- **Developing Distributed Smart Video Cameras (Phase II):**



- **Large-Scale Data Mining, Community Modeling, Human Behavior Modeling**
- *Also Course: EE6886 Multimedia Security Systems*

## Acknowledgements

- Navneet Panda (UCSB), Xiaodan Song (UW), Victor Sutan (CU), Jason Cardillo (CU)
- Olivier Verscheure (IBM), Lisa Amini (IBM)

Questions?

Contact: [cylin@ee.columbia.edu](mailto:cylin@ee.columbia.edu)  
<http://www.research.ibm.com/people/c/cylin>