# Cross-Modal Deep Networks for Document Image Classification

Souhail Bakkali[1]   Zuheng Ming[1]   Mickaël Coustaty[1]   Marçal Rusiñol[2]

[1]L3i, University of La Rochelle, France

[2]CVC, Universitat Autònoma de Barcelona, Spain
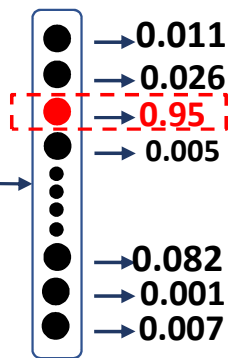
ICIP 2020

# Document Image Classification
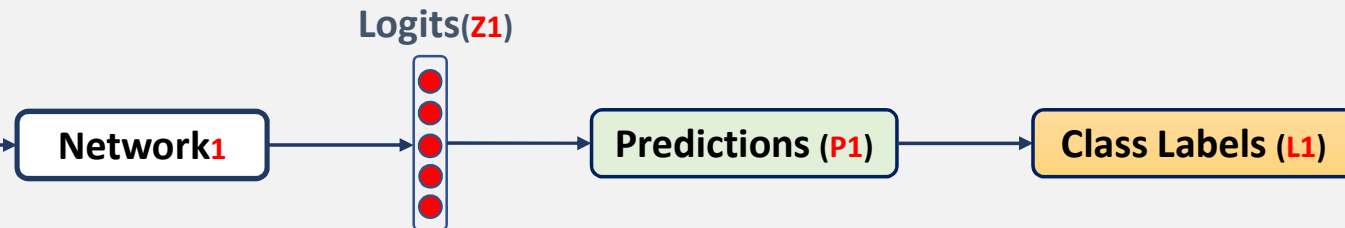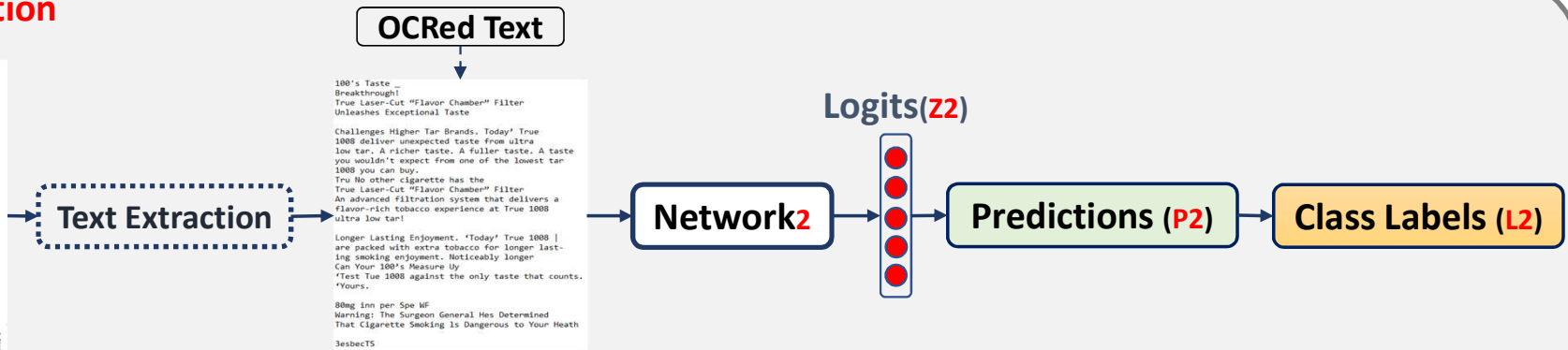


Input Document
(Advertisement)

Classifier

Class labels

0.011
0.026
0.95 → Advertisement
0.005

0.082
0.001
0.007

# Document Image Classification



**Image Classification**

Input Document → Network₁ → Logits(Z1) → Predictions (P1) → Class Labels (L1)

**Text Classification**

Input Document → Text Extraction → OCRed Text → Network₂ → Logits(Z2) → Predictions (P2) → Class Labels (L2)

# Motivation: **The Need for a Multimodal Analysis**

❑ Challenges:

- Some easy (stable) classes with similar visual content.

- Some hard classes with high variability.

- **10** classes with different complexity levels of the **Tobacco-3482** dataset.

❑ Solution:

- Need for a Multimodal Analysis



(a) Form

(b) Report

(c) News

(d) Email

# Architecture Network



**Image Stream**

100's Taste Breakthrough!

True Laser-Cut "Flavor Chamber" Filter
Unleashes Exceptional Taste

Visual Feature Extraction

Deep CNN

Text Extraction

**Text Stream**

Learning semantic information

Word Embedding Mechanism

OCRed Text

Cross-Modal Deep Network **(Fusion)** ❓

Logits

Softmax Cross-Entropy

# Image and Text Feature Learning

## Image Stream

**Visual Feature Extraction** ⇢ **Deep CNN**

### NasNet_Large

Conv2D · BatchNorm · Conv2D · BatchNorm · Average Pooling · BatchNorm · Conv2D · BatchNorm · MaxPool · Bottleneck · **Logits**

96 96 42 · 165 · 42 · 168 · 42 42 21 · 336 · 21 · 336 · 11 · 4032 · **d1**

## Text Stream

**Learning semantic information** → **Word Embedding Mechanism**

### Bert_base

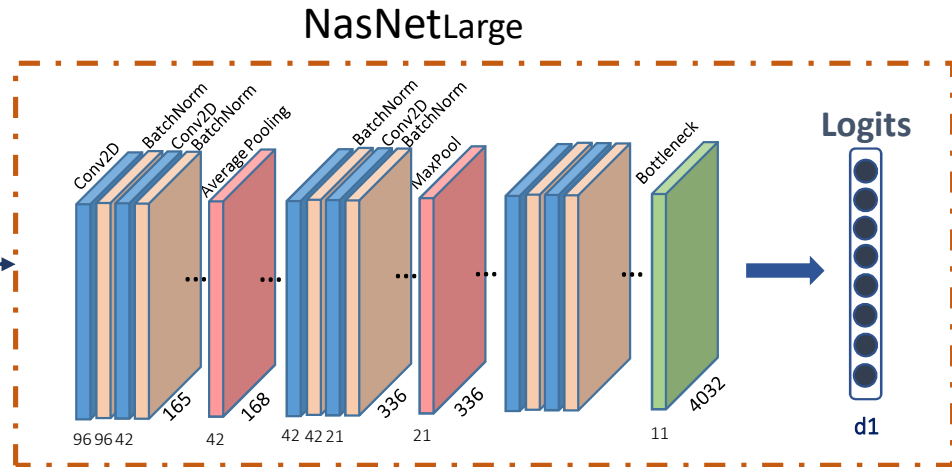| Input | [CLS] | Young | People | Form | Attitud | About | Smokin | By | Mid | Teens | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|

| Token Embeddings | $E_{[CLS]}$ | $E_{Young}$ | $E_{People}$ | $E_{Form}$ | $E_{Attitud}$ | $E_{About}$ | $E_{Smokin}$ | $E_{By}$ | $E_{Mid}$ | $E_{Teens}$ | $E_{[SEP]}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | + | + | + | + | + | + | + | + | + | + | + |
| Sentence Embedding | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |

**Logits**

**d2**

# How To Bridge These Two Different Modalities ?

# Cross-Modal Feature Learning

❑ **Naive Concatenation**



**Input Document**

NasNet$_{Large}$

Bert$_{Base}$

Logits(**Z1**)

d1

Logits(**Z2**)

d2

Fusion $\oplus$

Logits(**Z**)

d1 + d2

Class labels

10

The generated cross-modal features are given by:

$$X_a = [X_1, X_2], \quad X_a \in \mathbb{R}^{d_1 + d_2}$$

# Cross-Modal Feature Learning

❑ **Equal Concatenation**



Logits(**Z1**)

NasNet<sub>Large</sub>

d1

Logits(**Z**)

Class labels

Fusion ⊕

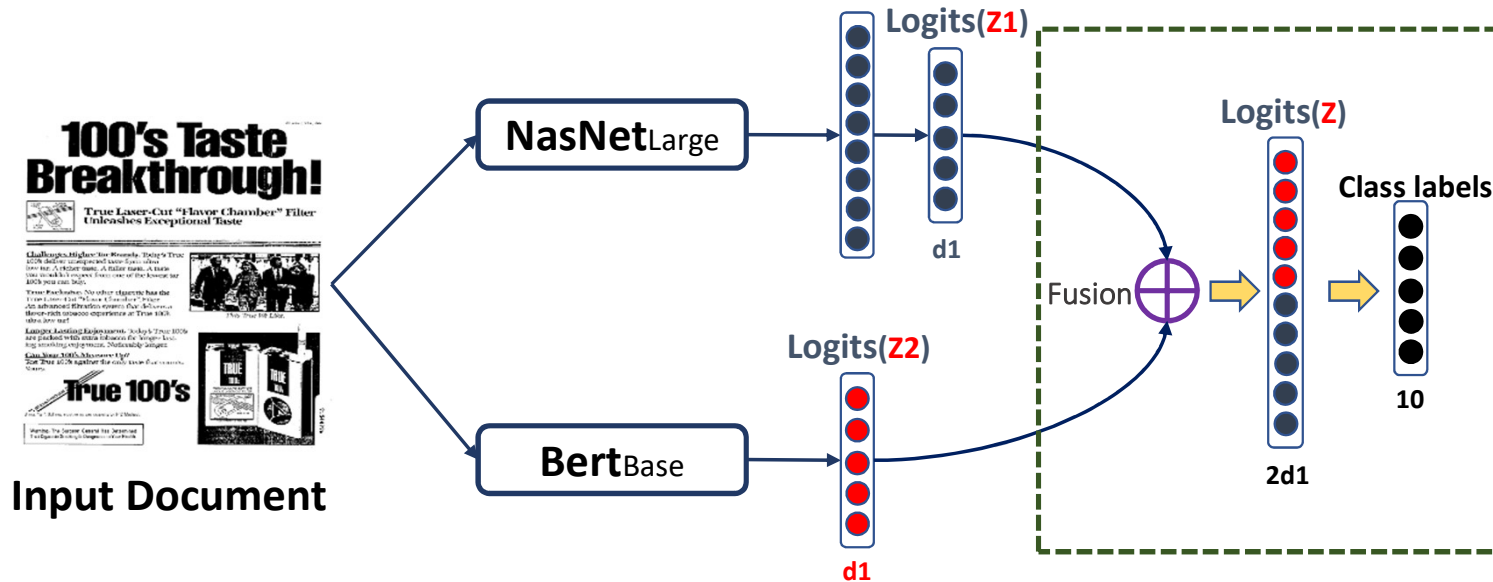2d1

10

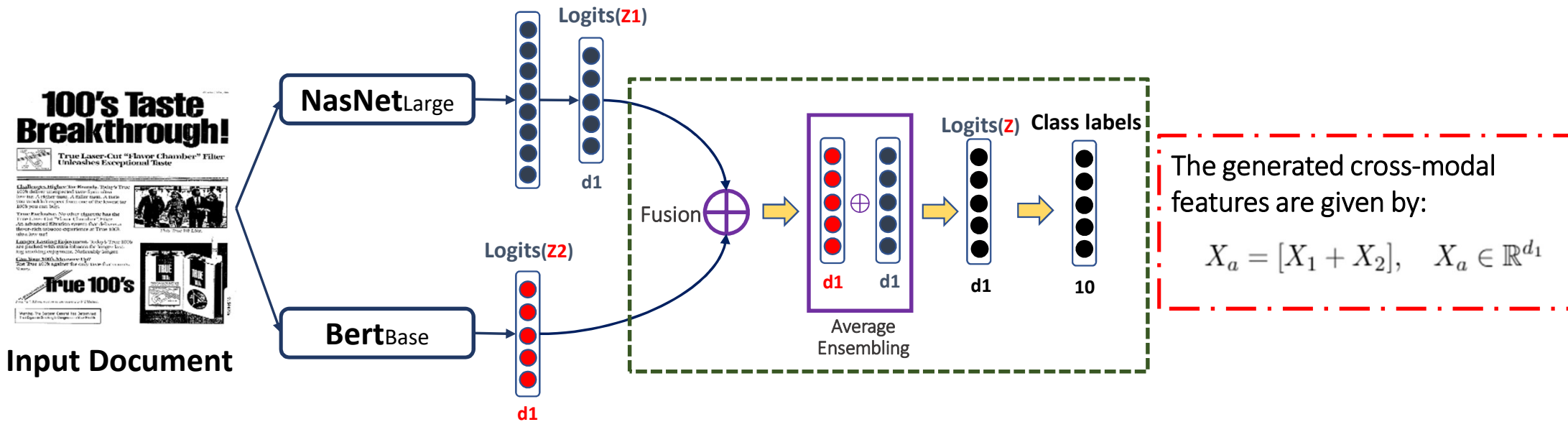Logits(**Z2**)

Bert<sub>Base</sub>

d1

**Input Document**

The generated cross-modal features are given by:

$$X_a = [X_1, X_2], \quad X_a \in \mathbb{R}^{2d_1}$$

# Cross-Modal Feature Learning

❑ **Superposing Concatenation**



**Logits(Z1)**

**NasNet**Large

**Input Document**

**Bert**Base

**Logits(Z2)**

d1

Fusion

d1    d1

Average
Ensembling

**Logits(Z)**    **Class labels**

d1    10

The generated cross-modal
features are given by:

$$X_a = [X_1 + X_2], \quad X_a \in \mathbb{R}^{d_1}$$

# Evaluation

**Table 1.** The Image-stream evaluation of best models from different methods on Tobacco-3482 dataset

| Method | Accuracy(%) |
|---|---|
| AlexNet [12] | 90.04 |
| GooGleNet [12] | 88.4 |
| VGG-16 [12] | 91.01 |
| ResNet-50 [12] | 91.13 |
| MobileNetV2 [9] | 84.50 |
| InceptionV3 [8] | 93.2 |
| **NASNet-Large** | 96.25 |

**Table 2.** Accuracy comparison of Text-stream state-of-the-art models on Tobacco-3482 dataset

| Method | Accuracy(%) |
|---|---|
| FastText-CNN [9] | 73.8 |
| Feature Ranking (ACC2) [8] | 87.1 |
| Glove-CNN1D-LSTM | 51 |
| Glove-GRU | 61 |
| **Bert** | 97.18 |

**Table 3.** Overall accuracy on the Tobacco-3482 dataset

| Model | Accuracy(%) | ADVE | Email | Form | Letter | Memo | News | Notes | Report | Resume | Scientific |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Single-Modal (Image-NASNet)** | 96.25 | 1 | 1 | 0.96 | 0.94 | 0.98 | 1 | 0.90 | 1 | 0.78 | 0.90 |
| **Single-Modal (Text-Bert)** | 97.18 | 0.97 | 0.99 | 0.98 | 0.93 | 0.97 | 0.98 | 0.89 | 1 | 0.96 | 0.95 |
| Multimodal Model [9] | 87.8 | 0.93 | 0.98 | 0.88 | 0.86 | 0.90 | 0.90 | 0.85 | 0.71 | 0.96 | 0.68 |
| Two Stream Model [8] | 95.8 | 0.94 | 0.98 | 0.95 | 0.98 | 0.97 | 0.97 | 0.88 | 0.92 | 1 | 0.93 |
| **Cross-Modal (Naive Concat.)** | 99.14 | 1 | 0.99 | 0.96 | 1 | 1 | 1 | 1 | 0.98 | 1 | 0.98 |
| **Cross-Modal (Equal Concat.)** | 98.42 | 0.98 | 0.99 | 0.95 | 1 | 0.98 | 0.97 | 1 | 1 | 0.96 | 0.98 |
| **Cross-Modal (Superposing fusion)** | 99.71 | 1 | 1 | 0.97 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

# Summary

- **Cross-Modal Network:** Learns simultaneously image and text features extracted from document images.

- Three **Feature Fusion** methods to perform Cross-Modal document image classification.

- **State-of-the-art** results compared to single-modal and multi-modal networks.