

# **Technology, Big Data, and Mixed Methods**

**Merlin Chowkwanyun**

**Big Data**



**Computing Horsepower**



**Volume**



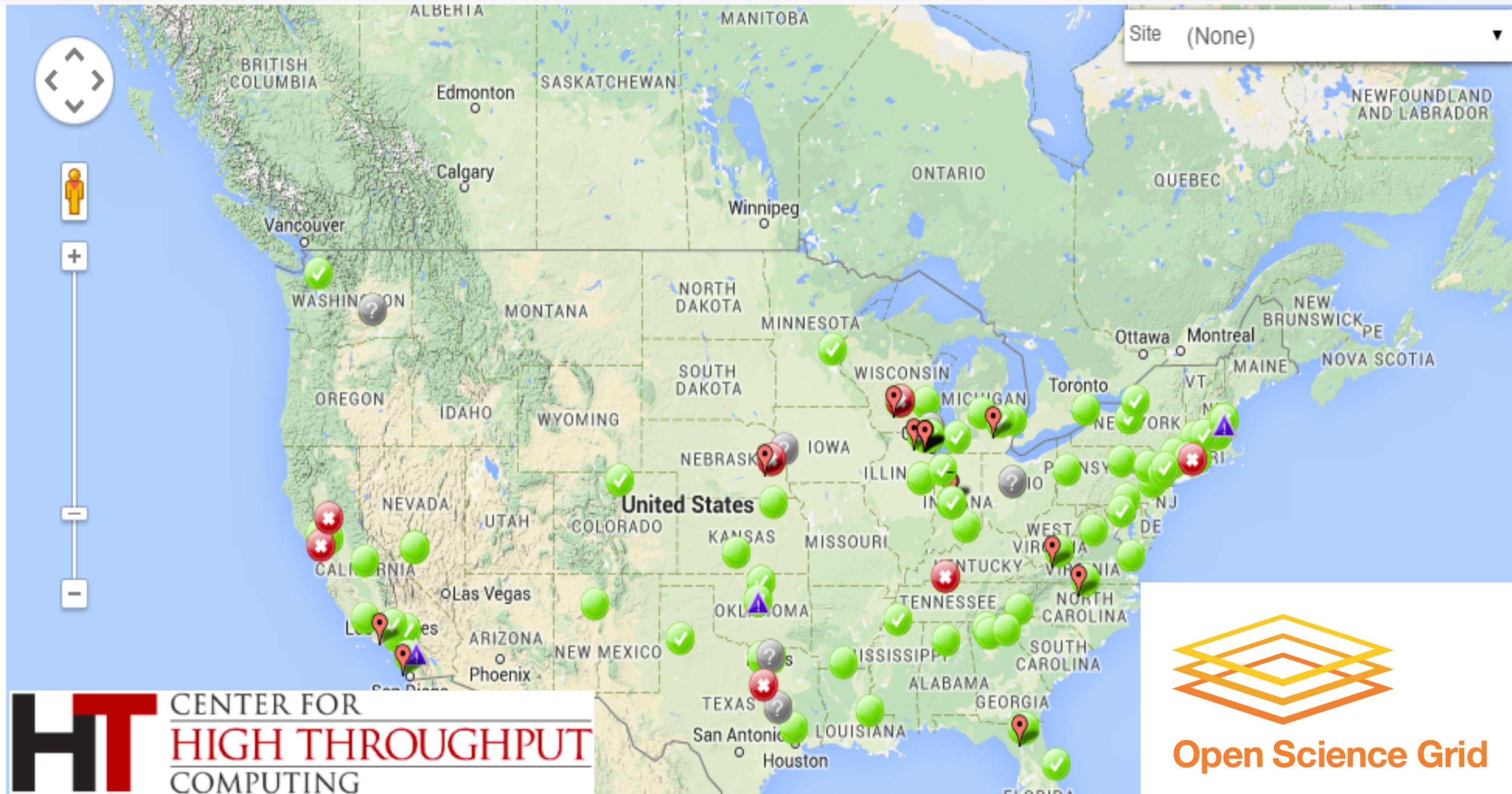
Site (None) ▾



CENTER FOR  
**HIGH THROUGHPUT**  
COMPUTING



Open Science Grid



Launch Instance

Connect

Actions ▾

Filter by tags and attributes or search by keyword

<input type="checkbox"/>	Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks	Alarm Status	Public DNS (IPv4)	IPv4 Public IP	IPv6 IPs	
<input type="checkbox"/>	GwenNERandNetwork	i-e917dc32	t2.medium	us-west-2b	<span style="color: green;">●</span> running	<span style="color: green;">✔</span> 2/2 checks ...	None		ec2-52-41-71-181.us-w...	52.41.71.181	-
<input type="checkbox"/>	webserver	i-14a7f7e2	c4.2xlarge	us-west-2b	<span style="color: green;">●</span> running	<span style="color: green;">✔</span> 2/2 checks ...	None		ec2-52-24-91-45.us-we...	52.24.91.45	-
<input type="checkbox"/>	WeekendProject1	i-0f3eb36e2375f5447	t2.small	us-west-2c	<span style="color: green;">●</span> running	<span style="color: green;">✔</span> 2/2 checks ...	None		ec2-34-208-188-226.us...	34.208.188.226	-
<input type="checkbox"/>	WeekendProject2	i-03330097cc491bc3e	t2.micro	us-west-2c	<span style="color: green;">●</span> running	Initializing	None		ec2-35-167-26-8.us-we...	35.167.26.8	-
<input type="checkbox"/>	WeekendProject3	i-089a483d1501dce9a	t2.micro	us-west-2c	<span style="color: green;">●</span> running	Initializing	None		ec2-34-209-146-144.us...	34.209.146.144	-
<input type="checkbox"/>	WeekendProject4	i-0d1e0ecbe3f7fe819	t2.micro	us-west-2c	<span style="color: green;">●</span> running	Initializing	None		ec2-54-148-12-199.us-w..	54.148.12.199	-
<input checked="" type="checkbox"/>	WeekendProject5	i-0ecca42fc8203798e	t2.micro	us-west-2c	<span style="color: green;">●</span> running	Initializing	None		ec2-35-163-251-89.us-w..	35.163.251.89	-

**Big Data**



**Computing Horsepower**

**Big Data**



**Computing Horsepower**



**Volume**

ID	Characteristic A	Characteristic B	Characteristic C	Characteristic D	Characteristic E	Characteristic F

# SQL

ID	Characteristic A	Characteristic B	Characteristic C	Characteristic D	Characteristic E	Characteristic F

# Non-Relational Flexible Databases

```
1  "retweet_count": 2301,  
2  "retweeted": false,  
3  "source": "<a href=\"http://twitter.com/download/android\" rel=\"nofollow\">Twitter for Android</a>",  
4  "text": "It is being reported by virtually everyone, and is a fact, that the media pile on against me is the worst in American political history!",  
5  "truncated": false,  
6  "user": {  
7    "contributors_enabled": false,  
8    "created_at": "Wed Mar 18 13:46:38 +0000 2009",  
9    "default_profile": false,  
10   "default_profile_image": false,  
11   "description": "#TrumpPence16",  
12   "entities": {  
13     "description": {  
14       "urls": []  
15     },  
16     "url": {  
17       "urls": [  
18         {  
19           "display_url": "DonaldJTrump.com",  
20           "expanded_url": "http://www.DonaldJTrump.com",  
21           "indices": [  
22             0,  
23             23  
24           ],  
25           "url": "https://t.co/mZB2hymxC9"  
26         }  
27       ]  
28     }  
29   },  
30   "favourites_count": 35,
```

# Non-Relational Flexible Databases

```
"retweet_count": 2301,
  "retweeted": false,
  "source": "<a href=\"http://twitter.com/download/android\" rel=\"nofollow\">Twitter for Android</a>",
  "text": "It is being reported by virtually everyone, and is a fact, that the media pile on against me is the worst in American political history!",
  "ideology": "conservative"
  "truncated": false,
  "user": {
    "contributors_enabled": false,
    "created_at": "Wed Mar 18 13:46:38 +0000 2009",
    "default_profile": false,
    "default_profile_image": false,
    "description": "#TrumpPence16",
    "entities": {
      "description": {
        "urls": []
      },
      "url": {
        "urls": [
          {
            "display_url": "DonaldJTrump.com",
            "expanded_url": "http://www.DonaldJTrump.com",
            "indices": [
              0,
              23
            ],
            "url": "https://t.co/mZB2hymxC9"
          }
        ]
      }
    }
  },
  "favourites_count": 35,
```

# INTRODUCING PROJECT TOXICDOCS ALPHA

Columbia University and the City University of New York

Millions of Previously Classified Documents on Industrial Poisons (and Counting)

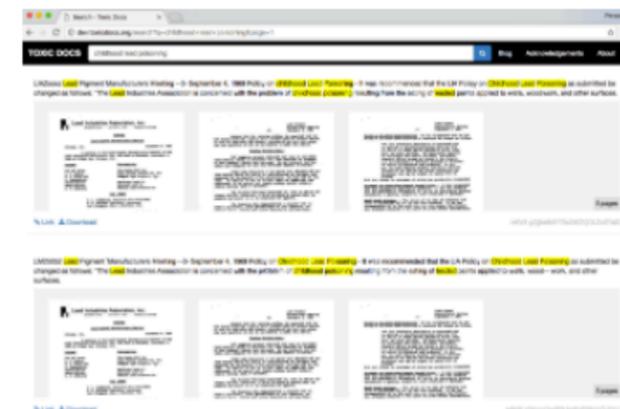


Advanced Search

## Blazingly fast searches of once-secret industry documents

Columbia University's Center for the History and Ethics of Public Health, located at its Mailman School of Public Health, and the City University of New York's Graduate Center are proud to jointly present Toxic Docs. This dataset and website contain millions of pages of previously secret documents about toxic substances. They include secret internal memoranda, emails, slides, board minutes, unpublished scientific studies, and expert witness reports -- among other kinds of documents -- that emerged in recent toxic tort litigation.

Over the next couple years, we'll be constantly adding material from lawsuits involving lead, asbestos, silica, and PCBs, among other dangerous substances. Innovations in parallel and cloud computing have made conversion of these documents into machine-readable, searchable text a far faster process than would have been the case just a decade ago.



# The Problem of Volume

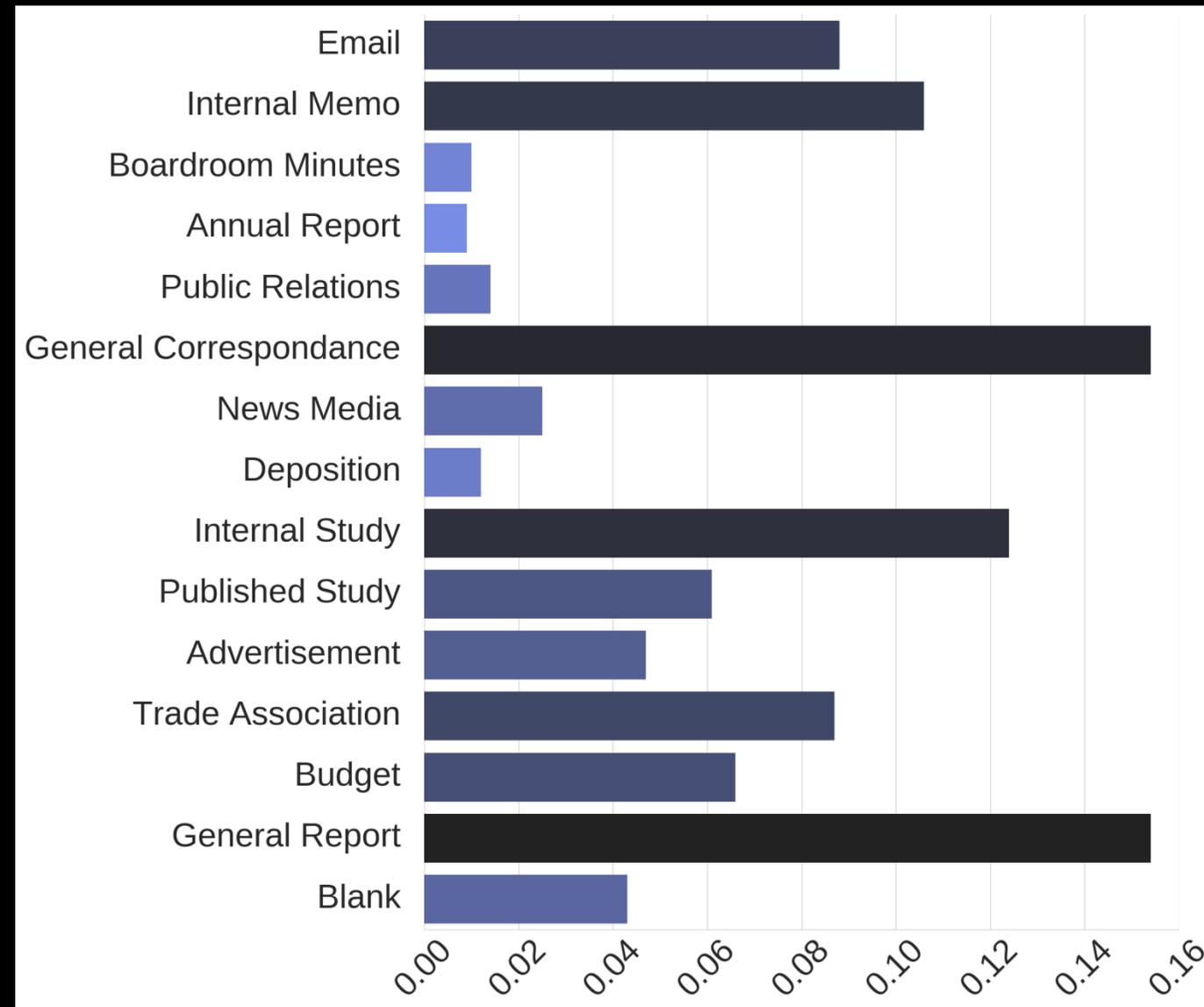
**1 computer**  
**4 million pages**  
**=**  
**6 months**

**4,228 servers**  
**4 million pages**  
**=**  
**1 day**

**Huge document  
stash**

**1k manually  
labeled**

**15 “genres”/categories**



"AGENDA"

# Board Minutes Example

"MINUTES"

AGENDA  
REGIONAL MEETING OF THE MCA BOARD OF DIRECTORS  
9:30 a.m. (PST), Wednesday, March 15, 1972  
Del Monte Lodge, Pebble Beach, California

- I. OPENING REMARKS AND INTRODUCTION OF GUESTS.
- II. MINUTES OF FEBRUARY 8, 1972, MEETING.
- III. FINANCIAL REPORT. (Enclosure)
- IV. BOARD OF DIRECTORS.
  - (a) Appointment of Alternate: C. A. Cash, Diamond Shamrock Corporation (for James A. Hughes)
  - (b) Proposed Rules of Organization and Procedure for the Economic Policy Review Committee. (Enclosure)
  - (c) Proposed Amendments to the Association's Bylaws. (Enclosure)
  - (d) Reports of Board Liaison Committees:  
Environmental Quality and Control -- Carl A. Gerstacker  
Government Relations -- George W. Russell  
Other Technical Committees -- Lee V. Dauler  
Public Relations and Education -- Luther S. Roehm
- V. COMMITTEE APPOINTMENTS. (Enclosure)
- VI. STAFF REPORT. (Appendix Enclosed)
- VII. ADJOURNMENT.

*See  
this  
shows  
membership*      *CHK;*

Next Meeting of the Directors -- The Madison, Washington, D. C., 8:30 a.m., (EST), Tuesday, April 11, 1972 (Breakfast Meeting).

CMA 036174

CONFIDENTIAL  
Subject to Protective Order in  
Ross v. Conoco, Inc., No. 90-4837      1386  
14th Judicial District Court  
Calcasieu Parish, Louisiana

MINUTES of the two hundred seventeenth meeting of the Directors of the Manufacturing Chemists' Association, Inc., held at the Pinnacle Club, New York City, Monday, November 20, 1972, at 4:00 p.m. There were present:

- Directors:
- |                           |                   |
|---------------------------|-------------------|
| F. Perry Wilson, Chairman | John M. Henske    |
| *H. Harold Bible          | H. E. Hirschland  |
| Frederick L. Bissinger    | James A. Hughes   |
| Werner C. Brown           | Edward R. Kane    |
| *Harry W. Buchanan        | Philip F. Kirk    |
| Fletcher L. Byrom         | Lloyd G. Lillico  |
| C. C. Candee              | Robert H. Malott  |
| Herschel H. Cudd          | *Robert M. Morris |
| Lee V. Dauler             | Charles S. Munson |
| Edward J. Donley          | Peter C. Reilly   |
| William C. Douce          | Luther S. Roehm   |
| *William P. Drake         | George W. Russell |
| William J. Driver         | Harold E. Thayer  |
| Carl A. Gerstacker        | *Jesse Werner     |
| *James M. Gill            | James R. Carnes   |

- Alternates:
- Warren M. Anderson (for F. Perry Wilson)
  - \*Stanley H. Anonsen (for Harold E. Thayer)
  - C. A. Cash (for James A. Hughes)
  - \*David N. Clark (for Lee V. Dauler)
  - \*Frank J. Connor (for H. E. Hirschland)
  - John T. Connor (for Frederick L. Bissinger)
  - \*Philip B. Dalton (for Jesse Werner)
  - \*H. N. Fiaccone (for Luther S. Roehm)
  - \*Richard Fleming (for Edward J. Donley)
  - John L. Gillis (for H. Harold Bible)
  - \*H. C. Hollands (for Lloyd G. Lillico)
  - L. H. Johnstone (for William C. Douce)
  - Gordon Kiddoo (for Donald G. Stevens)
  - \*James McWhirter (for William P. Drake)
  - \*Robert L. Mitchell (for John W. Brooks)
  - Richard M. Morrow (for Herschel H. Cudd)
  - Earl C. Ray (for Philip F. Kirk)
  - Thomas E. Reilly, Jr. (for Peter C. Reilly)
  - Dickson L. Whitney (for C. C. Candee)
  - \*Richard N. Williams (for John M. Henske)

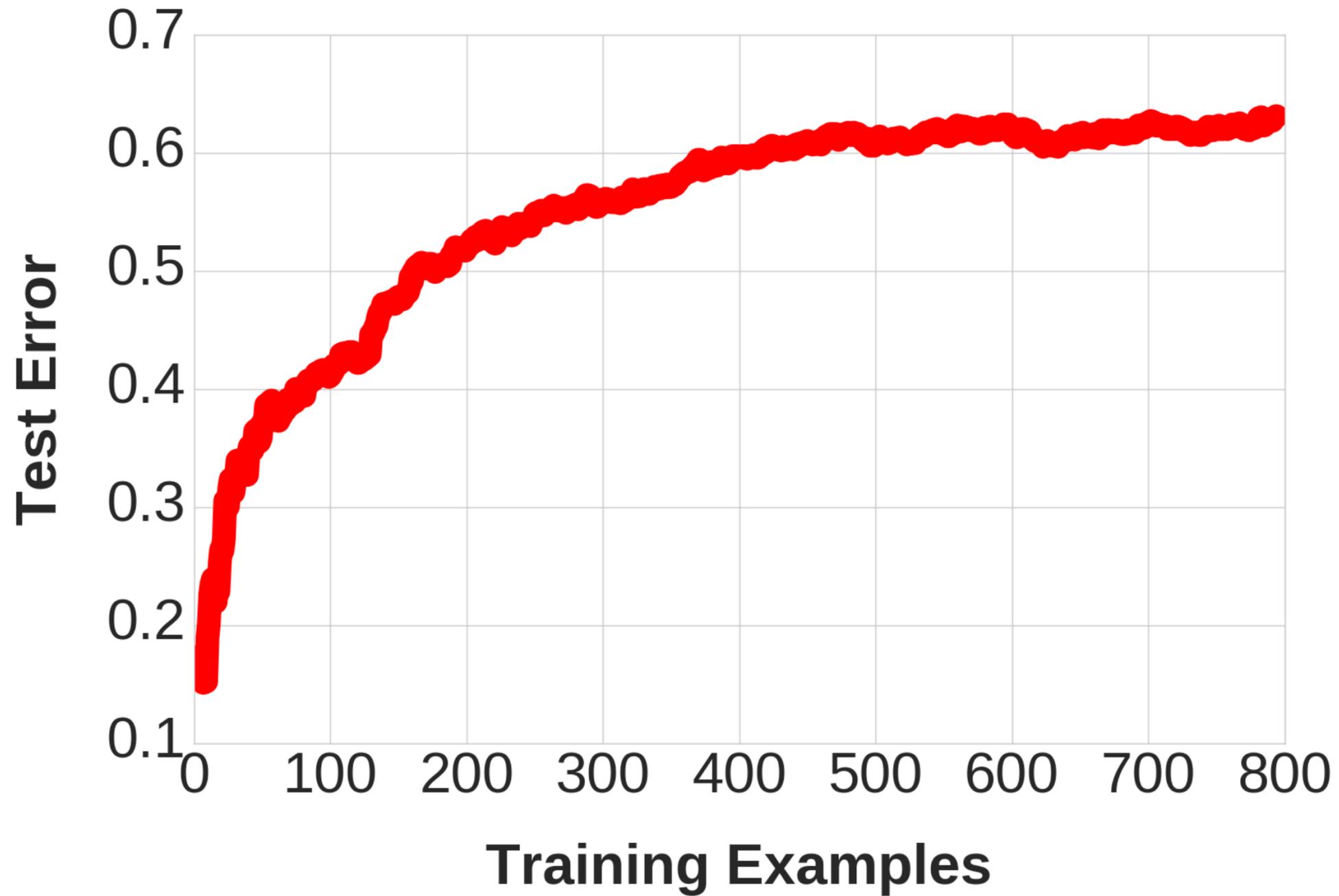
General Counsel: Lloyd Symington

By invitation: Bruce M. Barackman, MCA  
George E. Best, MCA

CMA 012526

Attendee Lists

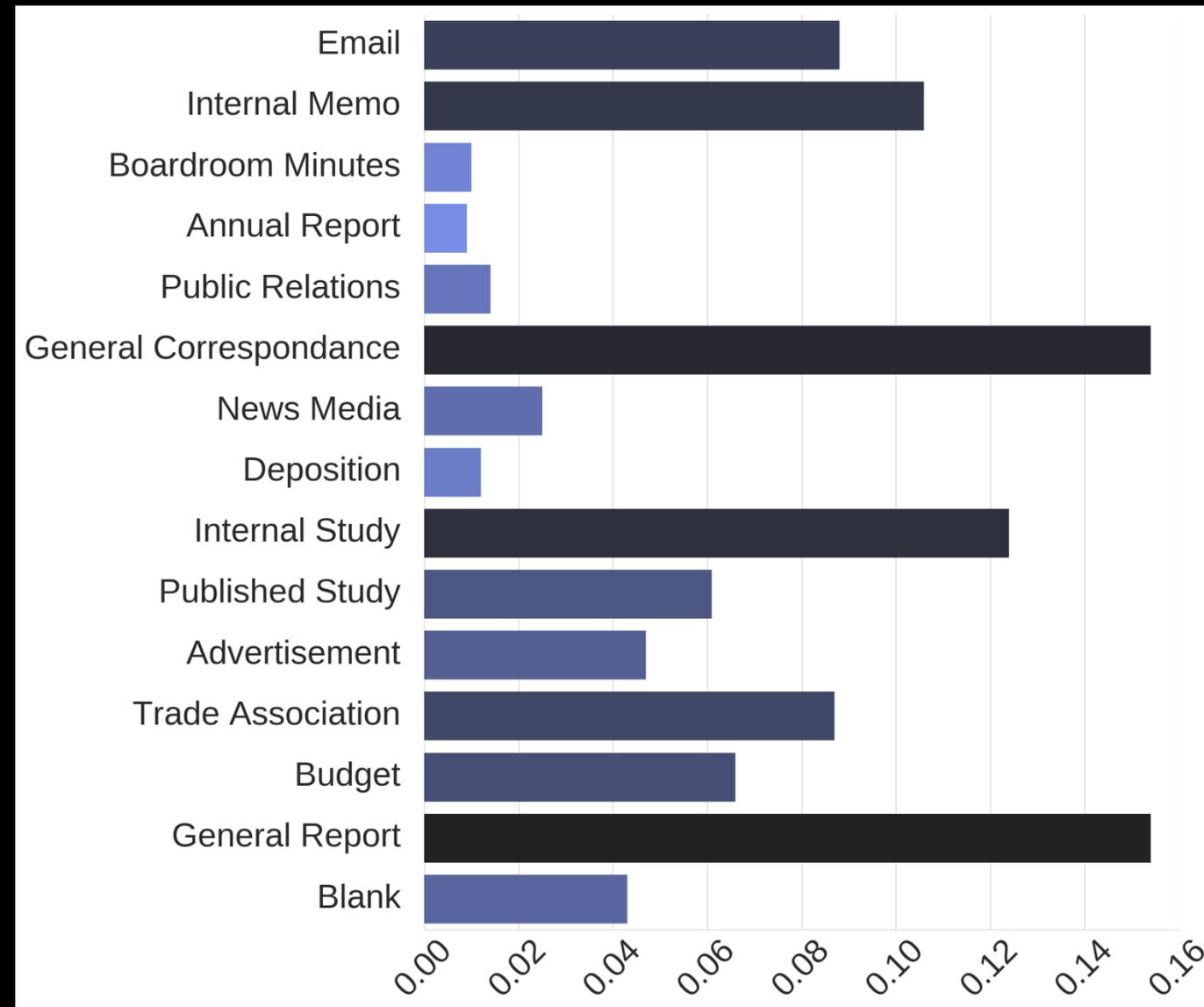
# SVM Learning Curve



**Huge document  
stash**

**1k manually  
labeled**

**15 “genres”/categories**





 new topic

Threads in Forum : The Foyer

Forum Tools

	Thread / Thread Starter	Rating	Last Post	Replies	Views
<b>Sticky Threads</b>					
	<a href="#">Let's not bicker and argue about who killed who, ...</a> Bernard		07-08-2015 09:30 AM by <a href="#">Bernard</a>	<a href="#">4</a>	15,359
	<a href="#">Read this before posting!</a> Bernard		03-19-2005 05:13 PM by <a href="#">Bernard</a>	<a href="#">0</a>	19,476
<b>Normal Threads</b>					
	<a href="#">Hello everyone!</a> Travis 71		Today 12:40 AM by <a href="#">Travis 71</a>	<a href="#">0</a>	1
	<a href="#">Wondering about keppra....</a> moldy73		04-26-2017 03:26 PM by <a href="#">Cint</a>	<a href="#">3</a>	87
	<a href="#">tongue spasms at night</a> whisper		04-26-2017 12:16 PM by <a href="#">whisper</a>	<a href="#">9</a>	634
	<a href="#">My history...</a> smitty1670		04-26-2017 12:14 PM by <a href="#">Sabbo</a>	<a href="#">4</a>	132
	<a href="#">Decision to concieve at 36 being diagnosed with Nocturnal Epilepsy....</a> Jnash23		04-26-2017 10:14 AM by <a href="#">Sabbo</a>	<a href="#">5</a>	119
	<a href="#">Hi! Everyone</a> CathyW		04-25-2017 12:10 AM by <a href="#">acshuman</a>	<a href="#">4</a>	94
	<a href="#">Meds Increase</a> smitty1670		04-24-2017 09:18 AM by <a href="#">Jeanbean</a>	<a href="#">5</a>	146

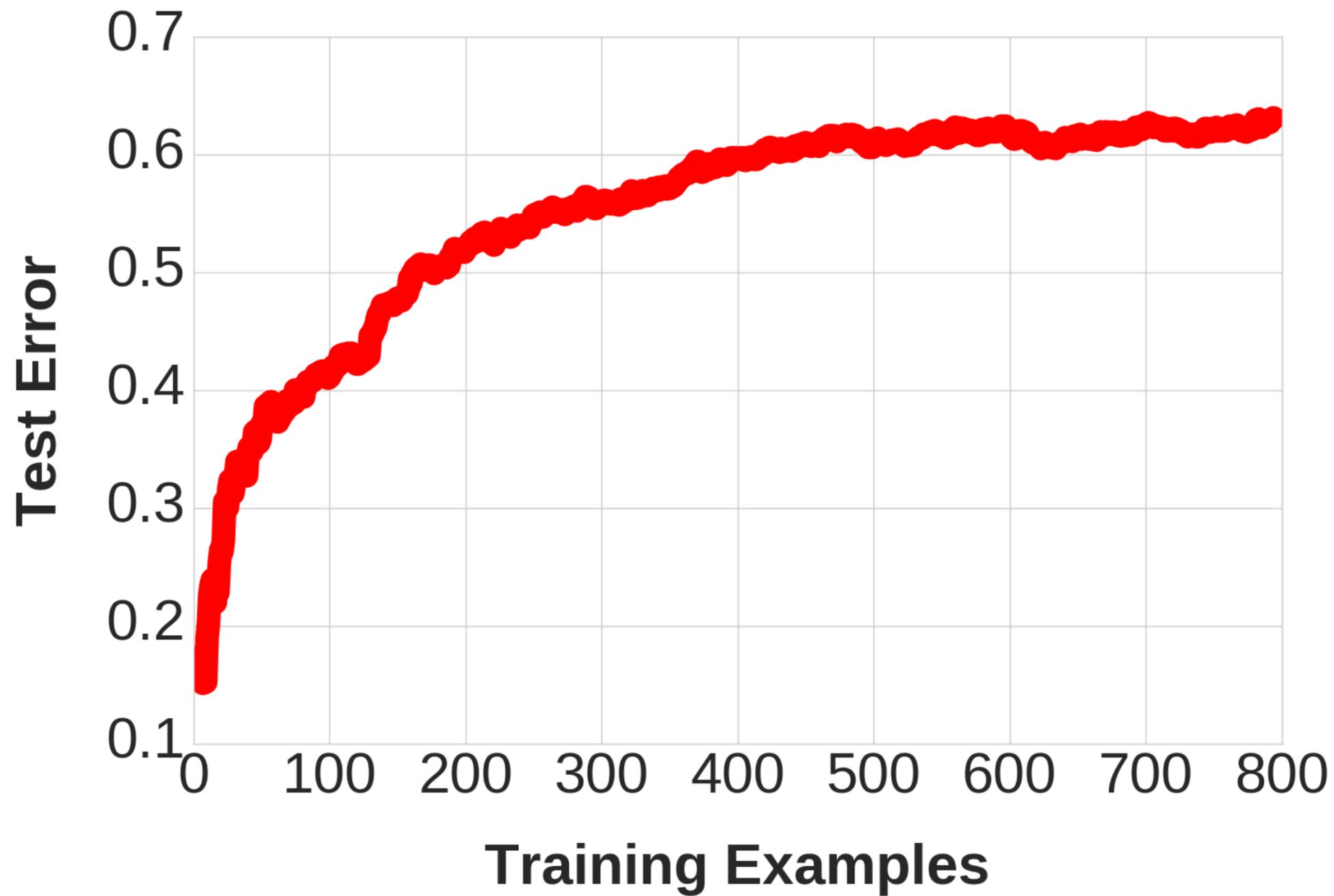


## List of Topics

1. jobs economy american ve businesses work tax americans year make ECONOMIC ISSUES
2. isil people syria military terrorist attacks united terrorists american ve ISIS/SYRIA
3. ve people make president question don obama lot things good JUNK
4. america people nation country americans history years page time american AMERICAN VALUES
5. ve work country make back working page day home great ECONOMIC ISSUES
6. iran nuclear international russia united world deal states weapons sanctions WEAPONS: RUSSIA/SYRIA
7. united states people countries world africa change democracy progress global DEMOCRACY ABROAD
8. war security america iraq afghanistan troops ve forces states support IRAQ/AFGHANISTAN
9. world people rights peace nations freedom future human israel copyright AMERICAN ALLIES
10. veterans service nation men american families honor serve lost military VETERANS
11. law information american court justice national process intelligence issues case JUNK
12. president obama people ve country political made page question continue JUNK

```
1 "contributors": null,
2   "coordinates": null,
3   "created_at": "Tue Aug 23 12:56:46 +0000 2016",
4   "entities": {
5     "hashtags": [],
6     "symbols": [],
7     "urls": [],
8     "user_mentions": [
9       {
10        "id": 1367531,
11        "id_str": "1367531",
12        "indices": [
13          72,
14          80
15        ],
16        "name": "Fox News",
17        "screen_name": "FoxNews"
18      }
19    ]
20  },
21  "favorite_count": 7826,
22  "favorited": false,
23  "geo": null,
24  "id": 768069472464666624,
25  "id_str": "768069472464666624",
26  "in_reply_to_screen_name": null,
27  "in_reply_to_status_id": null,
28  "in_reply_to_status_id_str": null,
29  "in_reply_to_user_id": null,
30  "in_reply_to_user_id_str": null,
31  "is_quote_status": false,
32  "lang": "en",
33  "place": null,
34  "retweet_count": 2547,
35  "retweeted": false,
36  "source": "<a href=\"http://twitter.com/download/android\" rel=\"nofollow\">Twitter for Android</a>",
37  "text": "I am now in Texas doing a big fundraiser for the Republican Party and a @FoxNews Special on the BORDER and with victims of border crime!",
38  "truncated": false,
39  "user": {
40    "contributors_enabled": false,
41    "created_at": "Wed Mar 18 13:46:38 +0000 2009",
42    "default_profile": false,
43    "default_profile_image": false,
44    "description": "#TrumpPence16",
45    "entities": {
46      "description": {
47        "urls": []
48      },
49      "url": {
50        "urls": [
51          {
52            "display_url": "DonaldJTrump.com",
53            "expanded_url": "http://www.DonaldJTrump.com",
54            "indices": [
55              0,
56              23
57            ],
58            "url": "https://t.co/mZB2hymxC9"
59          }
60        ]
61      }
62    }
63  }
64 }
```

# SVM Learning Curve





# OCR Accuracy

rate determination of a test compound during a cycle were as follows: A 20-1111. Sample of the mixed liquor (activated sludge + liquor) was withdrawn one hour after feeding and at the end of the aeration period. The sample of mixed liquor was then extracted and the extract concentrated according to the procedure given in Analytical Chemistry Method 71-18. The concentration of test compounds was then determined either by an ultraviolet (UV) spectroscopy or electron-capture gas chromatographic (EC-OC) method. Details of the UV methods are given in Analytical Chemistry Method 71-17 and the EC-OC methods in Analytical Chemistry Method 71-35. The disappearance rate was calculated from the following equation: