



# Conceptualizing Documents with Wikipedia

Tadashi Nomoto  
National Institute of Japanese  
Literature  
nomoto@acm.org

Noriko Kando  
National Institute of Informatics  
kando@nii.ac.jp

## Goal

The goal is to find a way to label a document cluster with a natural, informative phrase, together with some measure of confidence.

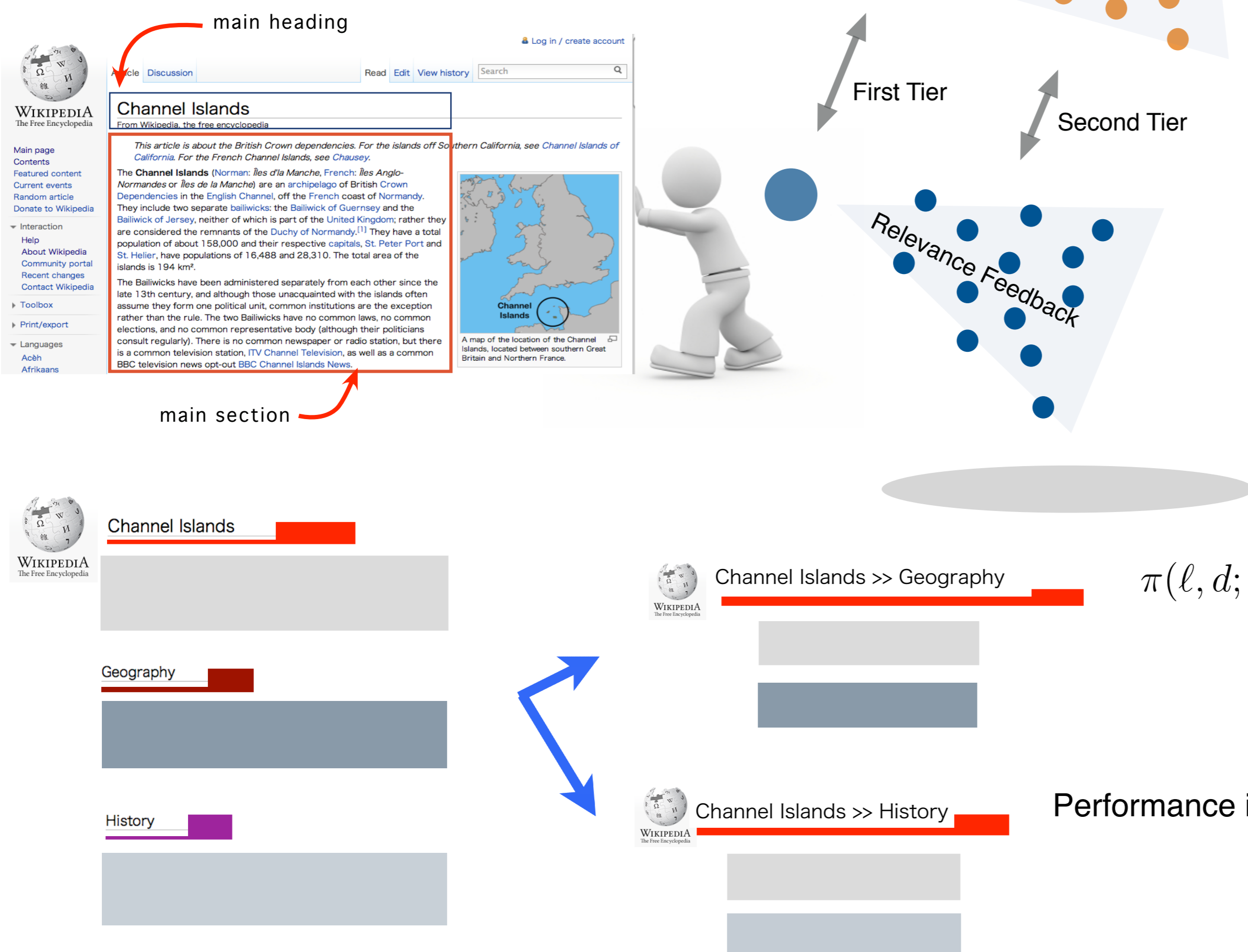


## Similarity Measures



COSINE	$C(\mathbf{q}, \mathbf{r}) = \frac{\mathbf{q}^T \mathbf{r}}{\ \mathbf{q}\ _2 \ \mathbf{r}\ _2}$
$L_1$ NORM ( $\ \mathbf{q} - \mathbf{r}\ _1$ )	$L_1(\mathbf{q}, \mathbf{r}) = \sum_t  \mathbf{q}(t) - \mathbf{r}(t) $
$L_2$ NORM ( $\ \mathbf{q} - \mathbf{r}\ _2^2$ )	$L_2(\mathbf{q}, \mathbf{r}) = \sum_t (\mathbf{q}(t) - \mathbf{r}(t))^2$
POLYNOMIAL KERNEL	$P(\mathbf{q}, \mathbf{r}) = (\mathbf{q}^T \mathbf{r} + C)^d$
RBF KERNEL	$R(\mathbf{q}, \mathbf{r}) = \exp(-\sigma^2 \ \mathbf{q} - \mathbf{r}\ _2^2)$
HELLINGER	$H(\mathbf{q}, \mathbf{r}) = \sum_t \sqrt{\mathbf{q}(t)\mathbf{r}(t)}$
SKREW DIVERGENCE	$Q(\mathbf{q}, \mathbf{r}) = D(\mathbf{q} \  \alpha \mathbf{q} + (1 - \alpha) \mathbf{r})$
SYMMETRIC KL	$S(\mathbf{q}, \mathbf{r}) = D(\mathbf{q} \  \mathbf{r}) + D(\mathbf{r} \  \mathbf{q})$

## Deconstructing Wikipedia



## Two-Tiered Similarity Model (TTM)

$$S(d_1, d_2) = \gamma \text{sim}_1(d_1, d_2) + (1 - \gamma) \text{sim}_2(\uparrow d_1, \uparrow d_2)$$

## Confidence Model (CFM)

$$P(R \geq t | \alpha_{\ell, d}) = \frac{P^*(\alpha_{\ell, d} | R \geq t) P^\dagger(R \geq t)}{P^*(\alpha_{\ell, d})}$$

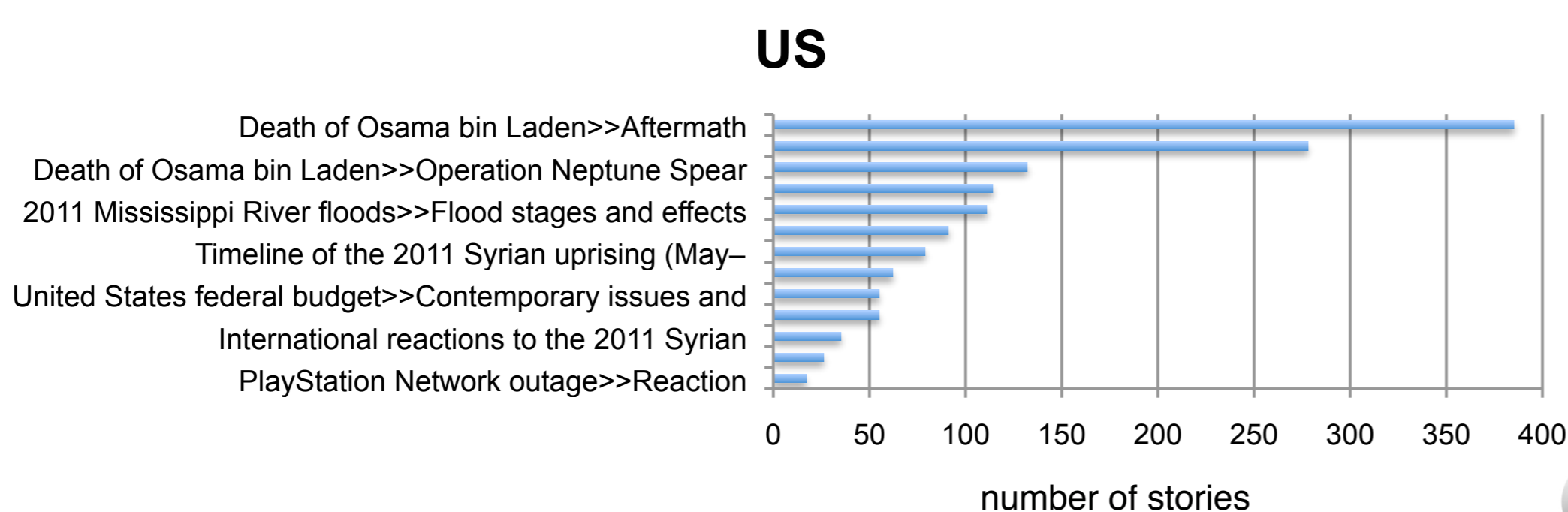
$$\pi(\ell, d; S) = \log P(R \geq t | S(g(\ell), d)) - \log P(R < t | S(g(\ell), d))$$

## Results

Performance in AUC with Wikipedia as a source for expansion at T2 (English)

MODEL	T1	T2	TTM	CFM(TTM)
$C$	0.678	0.665	0.693	<b>0.712</b>
$L_1$	0.684	0.664	0.689	<b>0.696</b>
$L_2$	0.673	0.694	0.707	<b>0.714</b>
$P$	<b>0.682</b>	0.612	0.632	0.677
$H$	0.682	0.612	0.631	<b>0.693</b>
$Q$	0.682	0.677	0.696	<b>0.709</b>
$R$	0.551	0.561	0.563	<b>0.633</b>
$S$	0.684	0.676	<b>0.689</b>	0.650
$\alpha$	1.0	0.5	0.0	MAX
ELN	0.667	0.671	0.679	0.679

## WikiLabel



$$l_\theta^* = \arg \max_{l: p[l] \in \mathcal{W}^M} \text{Score}(p[l], \theta|_N)$$

$$\text{Score}(p[l], \theta|_N) = \lambda \text{COS}(p[l], \theta|_N) + (1 - \lambda) \text{OVL}(l, \theta|_N)$$

