

# Introduction of MSRA Named Entity Corpus

## 1 Named Entity (NE) Tag-Set

Five NE categories with total 30 subcategories have been designed in the MSRA NE tag-set (see Table 1 below). The detail definition of each subcategory could be found in the spec, Tokenization Guidelines of Chinese Text.

Category	Subcategory	Tag-set of Format-1	Tag-set of Format-2
NAMEX	Person	P	PERSON
	Location	L	LOCATION
	Organization	O	ORGANIZATION
TIMEX	Date	dat	DATE
	Duration	dur	DURATION
	Time	tim	TIME
NUMEX	Percent	per	PERCENT
	Money	mon	MONEY
	Frequency	fre	FREQUENCY
	Integer	int	INTEGER
	Fraction	fra	FRACTION
	Decimal	dec	DECIMAL
	Ordinal	ord	ORDINAL
	Rate	rat	RATE
MEASUREX	Age	age	AGE
	Weight	wei	WEIGHT
	Length	len	LENGTH
	Temperature	tem	TEMPERATURE
	Angle	ang	ANGLE
	Area	are	AREA
	Capacity	cap	CAPACITY
	Speed	spe	SPEED
	Acceleration	acc	ACCELERATION
	Other measures	mea	MEASURE
ADDREX	Email	ema	EMAIL
	Phone	pho	PHONE
	Fax	fax	FAX
	Telex	tel	TELEX
	WWW	www	WWW
	Postalcode	pos	POSTALCODE

Table 1 Tag-set of Named Entities

## 2 Annotated Examples

While the format-1 is annotator oriented, the format-2 is XML-based which is the annotated format of the training and test corpora for the 3rd SAGHAN Bakeoff. Following are some annotated examples in these two formats respectively.

### [Example-1]

Format-1:

[dat 6月29日]、[dat 30日] [tim 晚上]，[L 北京市]下了[int 两场]大雨，笔者居住的宿舍楼前，宽[len 六七米]、长[len 30多米]的路上，积水达膝盖之上。

Format-2:

```
<sentence>
<w><TIMEX TYPE="DATE"> 6月29日</TIMEX></w><w>、</w><w><TIMEX
TYPE="DATE"> 30日</TIMEX></w><w><TIMEX TYPE="TIME"> 晚上
</TIMEX></w><w>,</w><w><NAMEX TYPE="LOCATION">北京市</NAMEX></w><w>
下</w><w>了</w><w><NUMEX TYPE="INTEGER">两场</NUMEX></w><w>大雨
</w><w>,</w><w>笔者</w><w>居住</w><w>的</w><w>宿舍</w><w>楼</w><w>前
</w><w>,</w><w>宽</w><w><MEASUREX TYPE="LENGTH">六七米
</MEASUREX></w><w>、</w><w>长</w><w><MEASUREX TYPE="LENGTH">30多米
</MEASUREX></w><w>的</w><w>路</w><w>上</w><w>,</w><w>积水</w><w>达
</w><w>膝盖</w><w>之上</w><w>。</w>
</sentence>
```

### [Example-2]

[dat 6月中下旬]，笔者到[L 意大利]、[L 西班牙]等国访问时，一个很深的感受是[L 意]、[L 西]两国的高速公路非常发达，东西南北，纵横成网，四通八达。

```
<sentence>
<w><TIMEX TYPE="DATE"> 6月中下旬</TIMEX></w><w>,</w><w>笔者</w><w>到
</w><w><NAMEX TYPE="LOCATION">意大利</NAMEX></w><w>、</w><w><NAMEX
TYPE="LOCATION">西班牙</NAMEX></w><w>等</w><w>国</w><w>访问</w><w>时
</w><w>,</w><w>一个</w><w>很</w><w>深</w><w>的</w><w>感受</w><w>是
</w><w><NAMEX TYPE="LOCATION">意</NAMEX></w><w>、</w><w><NAMEX
TYPE="LOCATION">西</NAMEX></w><w>两国</w><w>的</w><w>高速公路</w><w>非
常</w><w>发达</w><w>,</w><w>东西南北</w><w>,</w><w>纵横</w><w>成</w><w>
网</w><w>,</w><w>四通八达</w><w>。</w>
</sentence>
```

### [Example-3]

近 [dur 四年] 来，公司 [int 两个] 文明建设 >成果显著，主要经济指标保持以 [per 30 %] 的速度连年递增，先后评为全国 [int 五百家] 最佳经济效益企业、[L 北京市] 经济 [int 百] 强企业、首都文明单位标兵。

<sentence>

<w>近</w><w><TIMEX TYPE="DURATION">四年</TIMEX></w><w>来</w><w>，  
</w><w>公司</w><w><NUMEX TYPE="INTEGER">两个</NUMEX></w><w>文明  
</w><w>建设</w><w>成果</w><w>显著</w><w>，</w><w>主要</w><w>经济</w><w>  
指标</w><w>保持</w><w>以</w><w><NUMEX TYPE="PERCENT">30 %  
</NUMEX></w><w>的</w><w>速度</w><w>连年</w><w>递增</w><w>，</w><w>先后  
</w><w>被</w><w>评为</w><w>全国</w><w><NUMEX TYPE="INTEGER">五百家  
</NUMEX></w><w>最佳</w><w>经济效益</w><w>企业</w><w>、</w><w><NAMEX  
TYPE="LOCATION">北京市</NAMEX></w><w>经济</w><w><NUMEX  
TYPE="INTEGER">百</NUMEX></w><w>强</w><w>企业</w><w>、</w><w>首都  
</w><w>文明</w><w>单位</w><w>标兵</w><w>。</w>

</sentence>

### [Example-4]

[P 邓小平] 同志 [fre 多次] 强调，搞现代化建设、经济建设必须有政治保证。

<sentence>

<w><NAMEX TYPE="PERSON">邓小平</NAMEX></w><w>同志</w><w><NUMEX  
TYPE="FREQUENCY">多次</NUMEX></w><w>强调</w><w>，</w><w>搞</w><w>现代  
化</w><w>建设</w><w>、</w><w>经济</w><w>建设</w><w>必须</w><w>有</w><w>  
政治</w><w>保证</w><w>。</w>

</sentence>

### [Example-5]

我们 [O 北京城建集团一公司] 是 [dat 1983 年] 由基建工程兵集体转业组建的。

<sentence>

<w>我们</w><w><NAMEX TYPE="ORGANIZATION">北京城建集团一公司  
</NAMEX></w><w>是</w><w><TIMEX TYPE="DATE">1983 年</TIMEX></w><w>由  
</w><w>基建</w><w>工程兵</w><w>集体</w><w>转业</w><w>组建</w><w>的  
</w><w>。</w>

</sentence>

### [Example-6]

领导干部喜欢喝酒，他手下有 [int 几个] 『 [P 刘伶] 』、『 [P 李白] 』，也就确定无疑了。

<sentence>

<w>领导干部</w><w>喜欢</w><w>喝酒</w><w>，</w><w>他</w><w>手下</w><w>有

</w><w><NUMEX TYPE="INTEGER"> 几 个 </NUMEX></w><w> 『 </w><w><NAMEX TYPE="PERSON">刘 伶 </NAMEX></w><w>』 </w><w>、 </w><w>『 </w><w><NAMEX TYPE="PERSON">李 白 </NAMEX></w><w>』 </w><w>， </w><w>也 </w><w>就 </w><w>确定 </w><w>无疑 </w><w>了 </w><w>。 </w>  
</sentence>

#### [Example-7]

只有 这样 ， 才 能 顺乎 潮流 ， 经 受 挑战 ， 实 现 [L 中华]民族 的 伟 大 复 兴 。

<sentence>

<w> 只有 </w><w> 这样 </w><w>， </w><w> 才 </w><w> 能 </w><w> 顺乎 </w><w> 潮流 </w><w>， </w><w> 经 受 </w><w> 挑战 </w><w>， </w><w> 实 现 </w><w><NAMEX TYPE="LOCATION"> 中 华 </NAMEX> 民 族 </w><w> 的 </w><w> 伟 大 </w><w> 复 兴 </w><w>。 </w>

</sentence>

#### [Example-8]

科学研究 证明 ， 人 类 在 最近 [dur 三十年] 所 获 得 的 知 识 约 等 于 过 去 [dur 两千年] 之 总 和 ， 而 未 来 若 干 年 内 科 技 和 知 识 还 会 在 许 多 领 域 出 现 更 为 惊 人 的 突 破 。

<sentence>

<w> 科 学 研 究 </w><w> 证 明 </w><w>， </w><w> 人 类 </w><w> 在 </w><w> 最 近 </w><w><TIMEX TYPE="DURATION"> 三 十 年 </TIMEX></w><w> 所 </w><w> 获 得 </w><w> 的 </w><w> 知 识 </w><w> 约 </w><w> 等 于 </w><w> 过 去 </w><w><TIMEX TYPE="DURATION"> 两 千 年 </TIMEX></w><w> 之 </w><w> 总 和 </w><w>， </w><w> 而 </w><w> 未 来 </w><w> 若 干 </w><w> 年 </w><w> 内 </w><w> 科 技 </w><w> 和 </w><w> 知 识 </w><w> 还 </w><w> 会 </w><w> 在 </w><w> 许 多 </w><w> 领 域 </w><w> 出 现 </w><w> 更 为 </w><w> 惊 人 </w><w> 的 </w><w> 突 破 </w><w>。 </w>

</sentence>

#### [Example-9]

[L 建 行] 个 人 住 房 贷 款 [mon 万 元] 利 息 对 比 表

<sentence>

<w><NAMEX TYPE="ORGANIZATION"> 建 行 </NAMEX></w><w> 个 人 </w><w> 住 房 </w><w> 贷 款 </w><w><NUMEX TYPE="MONEY"> 万 元 </NUMEX></w><w> 利 息 </w><w> 负 担 </w><w> 对 比 </w><w> 表 </w>

</sentence>

#### [Example-10]

在 本 届 “ [L 灵 山] 杯 ” 评 选 中 ， 有 [int 二 十] 人 获 一 等 奖 ， [int 三 十 四] 人 获 [ord 二 等] 奖 ， [ord 三 等] 奖 [int 二 十 七 名] 。

<sentence>

<w> 在 </w><w> 本 届 </w><w> “ </w><w><NAMEX TYPE="LOCATION"> 灵 山 </NAMEX></w><w> 杯 </w><w> ” </w><w> 评 选 </w><w> 中 </w><w>， </w><w> 有

</w><w><NUMEX TYPE="INTEGER">二十</NUMEX></w><w>人</w><w>获</w><w>一  
 等奖</w><w>， </w><w><NUMEX TYPE="INTEGER">三十四</NUMEX></w><w>人  
 </w><w>获</w><w><NUMEX TYPE="ORDINAL">二等</NUMEX></w><w>奖</w><w>，  
 </w><w><NUMEX TYPE="ORDINAL">三等</NUMEX></w><w>奖</w><w><NUMEX  
 TYPE="INTEGER">二十七名</NUMEX></w><w>。</w>  
 </sentence>

#### [Example-11]

[dur 二十年] 来 ， 以 经 济 建 设 为 中 心 的 思 想 比 较 牢 固 ， 这 是 一  
 个 很 大 的 进 步 。

<w><TIMEX TYPE="DURATION">二十年</TIMEX></w><w>来</w><w>， </w><w>以  
 </w><w>经济</w><w>建设</w><w>为</w><w>中心</w><w>的</w><w>思想</w><w>比  
 较</w><w>牢固</w><w>， </w><w>这是</w><w>一个</w><w>很</w><w>大</w><w>的  
 </w><w>进步</w><w>。</w>  
 </sentence>

#### [Example-12]

[dat 7月1日] ， [O 中国人民银行] [ord 第五次] 降低 金融 机构 存 贷款 利  
 率 ， 中 长 期 贷 款 利 率 降 幅 达 到 或 超 过 [per 2个百分点] ， 将 有 助  
 于 刺 激 投 资 需 求 ， 加 快 我 国 基 础 设 施 建 设 。

<w><TIMEX TYPE="DATE">7月1日</TIMEX></w><w>， </w><w><NAMEX  
 TYPE="ORGANIZATION">中 国 人 民 银 行</NAMEX></w><w><NUMEX  
 TYPE="ORDINAL">第五次</NUMEX></w><w>降低</w><w>金融</w><w>机构</w><w>  
 存</w><w>贷款</w><w>利率</w><w>， </w><w>中 长 期</w><w>贷款</w><w>利率  
 </w><w>降 幅</w><w>达 到</w><w>或</w><w>超 过</w><w><NUMEX  
 TYPE="PERCENT">2个百分点</NUMEX></w><w>， </w><w>将</w><w>有 助 于</w><w>  
 刺 激</w><w>投 资</w><w>需 求</w><w>， </w><w>加 快</w><w>我 国</w><w>基 础  
 </w><w>设 施</w><w>建 设</w><w>。</w>  
 </sentence>

#### [Example-13]

[O 北京大学光华管理学院] 博士生 [P >张后奇]

<w><NAMEX TYPE="ORGANIZATION">北京大学光华管理学院</NAMEX></w><w>博士  
 生</w><w><NAMEX TYPE="PERSON">张后奇</w><w>  
 </sentence>

#### [Example-14]

[L 中国建设银行房贷部] 总经理 [P 李庆振]

<w><NAMEX TYPE="ORGANIZATION">中国建设银行房贷部</NAMEX></w><w>总 经 理  
 </w><w><NAMEX TYPE="PERSON">李庆振</w><w>

</sentence>

#### [Example-15]

我们 一定 要 认真 做好 [dat 下半年] 的 工作 ， 确保 [O 首钢总公司 全年  
实现 利润 [mon 9.4 亿元] 的 目标 ， 并 >力争 多 完成 一些 。

<sentence>

<w> 我们 </w><w> 一定 </w><w> 要 </w><w> 认真 </w><w> 做好 </w><w><TIMEX  
TYPE="DATE"> 下半年 </TIMEX></w><w> 的 </w><w> 工作 </w><w> , </w><w> 确保  
</w><w><NAMEX TYPE="ORGANIZATION"> 首钢总公司 </NAMEX></w><w> 全年  
</w><w> 实现 </w><w> 利 润 </w><w><NUMEX TYPE="MONEY"> 9 . 4 亿 元  
</NUMEX></w><w> 的 </w><w> 目标 </w><w> , </w><w> 并 </w><w> 力争 </w><w> 多  
</w><w>完成</w><w>一些</w><w>。 </w>

</sentence>

#### [Example-16]

[O 县委] 决定 选派 任 了 [dur 八年] [O 城建局] 长 的 [P 周欣光] 担任  
[老干部局] 长 。

<sentence>

<w><NAMEX TYPE="ORGANIZATION"> 县委 </NAMEX></w><w> 决定 </w><w> 选派  
</w><w> 任 </w><w> 了 </w><w><TIMEX TYPE="DURATION"> 八 年  
</TIMEX></w><w><NAMEX TYPE="ORGANIZATION"> 城建局 </NAMEX></w><w> 长  
</w><w> 的 </w><w><NAMEX TYPE="PERSON"> 周 欣 光 </NAMEX></w><w> 担任  
</w><w><NAMEX TYPE="ORGANIZATION"> 老干部局 </NAMEX></w><w>长</w><w>。

</w>

</sentence>

#### [Example-17]

[L 喇嘛寺村]/ 地处 [L 承德避暑山庄] ， [L 山庄] 寺庙 林立 ， 僧侣 穿梭 ，  
[L 山庄] [L 外八庙] 的 [ord 第一个] 庙 就 是 [L 喇嘛寺] 。

<sentence>

<w><NAMEX TYPE="LOCATION">喇嘛寺村</NAMEX></w><w>地处</w><w><NAMEX  
TYPE="LOCATION"> 承 德 避 暑 山 庄 </NAMEX></w><w> , </w><w><NAMEX  
TYPE="LOCATION">山庄</NAMEX></w><w>寺庙</w><w>林立</w><w>, </w><w>僧侣  
</w><w> 穿 梭 </w><w> , </w><w><NAMEX TYPE="LOCATION"> 山 庄  
</NAMEX></w><w><NAMEX TYPE="LOCATION"> 外 八 庙 </NAMEX></w><w> 的  
</w><w><NUMEX TYPE="ORDINAL">第一个</NUMEX></w><w>庙</w><w>就</w><w>  
是</w><w><NAMEX TYPE="LOCATION">喇嘛寺</NAMEX></w><w>。 </w>

</sentence>

### 3 Statistics of Training Corpus

# of sentences	# of words (including punctuations)	# of NE tags	Ave. # of words per sentence	Ave. # of NE tags per sentence
46,364	1,265,764	118,643	27.30	2.56

Table 2 Size of MSRA Training Corpus

Number of L	Number of O	Number of P	Total number
36,860	20,584	17,615	75,059

Table 3 Number of NAMEX Tags

Number of dat	Number of dur	Number of tim	Total number
12,661	3,928	1,335	17,924

Table 4 Number of TIMEX Tags

# of int	# of ord	# of mon	# of per	# of fre
10,433	3,506	2,346	1,786	939
# of rat	# of fra	# of dec	Total number	
622	467	340	20,439	

Table 5 Number of NUMEX Tags

# of tem	# of len	# of age	# of are	# of wei	Total number
2,741	735	612	439	424	5,165
# of mea	# of cap	# of ang	# of spe	# of acc	
137	46	18	13	0	

Table 6 Number of MEASUREX Tags

# of pho	# of www	# of pos	# of fax	# of tel	Total number
40	13	3	0	0	56

Table 7 Number of ADDREX Tags

<END>