

## Supplementary Information

## 高质量的材料科学文本挖掘数据集构建方法\*

刘悦<sup>1)4)</sup> 刘大晖<sup>1)</sup> 葛献远<sup>1)</sup> 杨正伟<sup>1)</sup> 马舒畅<sup>1)</sup> 邹喆义<sup>5)</sup> 施思齐<sup>2)3)</sup>

1) (上海大学计算机工程与科学学院, 上海 200444)

2) (上海大学材料科学与工程学院, 上海 200444)

3) (上海大学材料基因组工程研究院, 上海 200444)

4) (上海市智能计算系统工程技术研究中心, 上海 200444)

5) (湘潭大学材料科学与工程学院, 湘潭 411105)

表 S1 NASICON 实体识别数据集的注释示例

Table S1. Examples of annotation for entity recognition dataset of NASICON.

实体类型	示例
Composition	The bond valence sum (BVS) of 1.08 is in good agreement with the value expected for <b>Na+</b> .
Structure	The results from both measurements can be correlated only, if at room temperature <b>Na+</b> is located in the interstitial sites and if the transition of these ions to the <b>Na2 site</b> occurs at high temperature.
Property	Conductivity measurements in the range 30 – 350 <sYm> reveal the <b>activation energy</b> of 0.3 eV for <b>Na+</b> conduction but <b>conductivity</b> values were found to change with temperature of sample preparation.
Processing	<b>Regrinding</b> and <b>reheating</b> the mix results in very slow incorporation of the free ZrO <sub>2</sub> , probably by replacement of <b>Na+</b> from the zirconium sites.
Application	It is shown that Na <sub>3</sub> TiP <sub>3</sub> O <sub>9</sub> N can reversibly cycle Na-ions in a manner suitable for <b>secondary batteries</b> , and that the volume changes on Na removal are remarkably small (<1%) with respect to other known Na-ion.
Characterization	The <b>XRD diffraction pattern</b> of Na <sub>3</sub> MnTi(PO <sub>4</sub> ) <sub>3</sub> can be indexed into a rhombohedral NASICON type unit cell with the R3c space group.
Feature	<b>Single crystal</b> x-ray analysis is used to identify the composition NaZr <sub>2</sub> P <sub>3</sub> O <sub>12</sub> and to refine its structure, which is rhombohedral.
Condition	The anisotropy of the thermal vibrations of sodium atoms in NaSn <sub>2</sub> (PO <sub>4</sub> ) <sub>3</sub> at <b>room temperature</b> is described by two different flattened ellipsoids.

表 S2 NASICON 关系抽取数据集的注释示例

**Table S2.** Example of annotation for relational extraction dataset of NASICON.

关系类型	示例
Cause-Effect	<p>1. (<i>Composition-Property, Property-Property</i>) In the system <math>\text{Na}_{(1+x)}\text{Zr}_2\text{SixP}_{(3-x)}\text{O}_{12}</math>, on the other hand, the introduction of excess <b>Na ions</b><sup>(1)</sup> causes electrostatic <b>Na<sup>+</sup>-Na<sup>+</sup> interactions</b><sup>(2)</sup> that can lower the <b>activation energy</b><sup>(3)</sup> even though transport must be via a NaI site.</p> <p>2. (<i>Structure-Property</i>) The increase of the <b>M1 site size</b><sup>(1)</sup> in <math>\text{Na}_2\text{SnFe}(\text{PO}_4)_3</math> is accompanied by oxygen displacements perpendicular to the c axis which gives rise to <b>rotation</b><sup>(2)</sup> of the <math>\text{PO}_4</math> tetrahedra and leads to a <b>distortion</b><sup>(3)</sup> of the <math>\text{Sn}(\text{Fe})(1-x)(\text{PO}_4)_3</math> framework.</p>
Component-Whole	<p>1. (<i>Composition-Composition</i>) The original Na super ionic conductors NASICON materials are solid solutions derived from <b>NaZr2P3O12</b><sup>(1)</sup> by partial replacement of <b>P</b><sup>(2)</sup> by <b>Si</b><sup>(3)</sup> with extra <b>Na</b><sup>(4)</sup> to balance the charges.</p>
Feature-Of	<p>1. (<i>Composition-Property</i>) The calculated BVS value 5.22 shows that the <b>As5 cation</b><sup>(1)</sup> is also slightly <b>over bonded</b><sup>(2)</sup>.</p> <p>2. (<i>Composition-Structure</i>) The present paper reports on the <b>crystal structure</b> and vibrational spectra of <b>NaZr2(AsO4)3</b>.</p>
Located-Of	<p>1. (<i>Composition-Structure</i>) The <b>Na cation</b><sup>(1)</sup> is located on the <b>6b position</b><sup>(2)</sup> with a trigonal antiprismatic coordination and enhanced anisotropic displacement parameters.</p>
Instance-Of	<p>1. (<i>Property-Property</i>) The monochromator is a crystal of Ge that selects a <b>wavelength</b><sup>(1)</sup> of <b>1.594 Å</b><sup>(2)</sup>.</p> <p>2. (<i>Feature-Feature</i>) X-ray powder diffraction shows that the <b>phosphates</b><sup>(1)</sup> belong to the <b>NZP type</b><sup>(2)</sup>.</p> <p>3. (<i>Structure-Structure</i>) The four <b>P-O bonds</b><sup>(1)</sup> in the near regular <b>tetrahedron</b><sup>(2)</sup> (point symmetry<sub>2</sub>) range from 1.524 to 1.525 Å, with O-P-O angles deviating by no more than 1.5 &lt;sym&gt; from the ideal 109.5 &lt;sym&gt; tetrahedral angle (Table III).</p>
Condition-On	<p>1. (<i>Processing-Condition</i>) For each crystal structure determination, data are collected using <b>Mo Ka radiation</b><sup>(1)</sup> up to <b>29-650</b><sup>(2)</sup>.</p> <p>2. (<i>Property-Condition</i>) The <b>disorder</b><sup>(1)</sup> is larger at <b>100 K</b><sup>(2)</sup> than at <b>295 K</b><sup>(3)</sup>.</p>
Method-Of	<p>1. (<i>Characterization-Property</i>) The <b>Rietveld plots</b><sup>(1)</sup> represent a good structure fit between observed and calculated intensity with satisfactory <b>R-factors</b><sup>(2)</sup>.</p> <p>2. (<i>Characterization-Property</i>) Its <b>anisotropic thermal expansion</b><sup>(1)</sup> has been calculated from high temperature <b>X-ray diffraction</b><sup>(2)</sup>, and it is linear in a range from room temperature up to 800 &lt;sym&gt;.</p> <p>3. (<i>Processing-Condition</i>) <b>Good crystals</b><sup>(1)</sup> can, however, be obtained after <b>tempering</b><sup>(2)</sup> in platinum crucible for several weeks at 11000 &lt;sym&gt;.</p>
Other	<p>1. (<i>Condition-Property</i>) At <b>room temperature</b><sup>(1)</sup> no <b>diffuse intensity</b><sup>(2)</sup> is observed.</p>

表 S3 MatBERT-BiLSTM-CRF 模型参数

**Table S3.** Parameters of MatBERT-BiLSTM-CRF.

参数名称	值
批量大小 (Batch size)	32
迭代周期 (Epoch)	100
单词向量维度 (Word vector dimension)	768
LSTM 单元维度 (LSTM unit dimension)	128
丢包率 (Dropout rate)	0.1
学习率 (Learning rate)	$3 \times 10^{-5}$
优化器 (Optimizer)	AdamW
提前停止耐力值 (Early stopping patience)	3
最大句子长度 (Max sentence length)	75

表 S4 不同机器学习模型的激活能预测结果<sup>[33]</sup>

**Table S4.** Results of different machine learning models<sup>[33]</sup>.

Models	Dataset <sub>31</sub>			Dataset <sub>45</sub>		
	RMSE	MAPE	$R^2$	RMSE	MAPE	$R^2$
LASSO	0.09	0.06	0.86	0.06	0.04	0.94
GPR	0.09	0.06	0.86	0.05	0.04	0.96
Ridge	0.09	0.06	0.86	0.05	0.04	0.95
SVR	0.10	0.07	0.84	0.07	0.06	0.92
KNN	0.11	0.07	0.80	0.09	0.06	0.87
RF	0.10	0.06	0.83	0.05	0.04	0.96