

# 数据集构建与使用说明

## 一、训练集（8 种字体，字号选择 2 种，小四和三号，每个字符 16 个样本）

### 1、宋体：

哈尔滨工业大学的校训是规格严格功夫到家

### 2、仿宋：

哈尔滨工业大学的校训是规格严格功夫到家

### 3、微软雅黑：

哈尔滨工业大学的校训是规格严格功夫到家

### 4、华光隶书：

哈尔滨工业大学的校训是规格严格功夫到家

### 5、华光姚体：

哈尔滨工业大学的校训是规格严格功夫到家

### 6、华光彩云：

哈尔滨工业大学的校训是规格严格功夫到家

### 7、华光行书：

哈尔滨工业大学的校训是规格严格功夫到家

### 8、华光行楷：

哈尔滨工业大学的校训是规格严格功夫到家

## 二、测试集（4 种字体，字号选择 2 种，小四和三号，每个字符 8 个样本）

使用了与训练集不同的另外 4 种字体，具体内容略。

## 三、数据集构建流程（已经提供，无需实现）：

- 1、基于训练集和测试集的内容要求，打印输出训练集和测试集文档样张，然后基于扫描仪或手机等方式转换为图片格式。
- 2、自动或手动切割图片中的文字，得到单个字符图像样本。

## 四、数据集使用

- 1、将 train\_small 和 train\_large 进行合并，构成完整训练数据集。根据样本顺序特点，构建每个字符图片的对应标签。
- 2、提供的字符样本图片是原始灰度图进行了二值化处理后的版本，质量稍有降低。
- 3、可选步骤：训练数据集可根据需要进行扩充，如增加不同的字号，通过扫描仪、手机等不同方式及不同的采集参数将电子文档或纸质文档转换为图片格式。如果采用了数据扩充，提交文档中必须对具体方式进行说明。
- 4、可选步骤：训练数据集也可进行一定的数据增强，如图像变换或增加噪声等。如果采用了数据增强，提交文档中必须对具体方式进行说明。
- 5、训练集样本进行扩充或增强时，不得使用测试集中的字体类型或测试集中的样本。
- 6、测试集样本不得做任何改动。

# 文字样本样张制作

## 训练集 1

Train\_small

### 文字版

哈尔滨工业大学的校训是规格严格功夫到家

哈尔滨工业大学的校训是规格严格功夫到家

哈尔滨工业大学的校训是规格严格功夫到家

哈尔滨工业大学的校训是规格严格功夫到家

哈尔滨工业大学的校训是规格严格功夫到家

哈尔滨工业大学的校训是规格严格功夫到家

哈尔滨工业大学的校训是规格严格功夫到家

哈尔滨工业大学的校训是规格严格功夫到家

### 图片版：（二值 8 位灰度图）



## 训练集 2

Train\_Large

### 文字版

哈尔滨工业大学的校训是规格严格功夫到家  
哈尔滨工业大学的校训是规格严格功夫到家  
哈尔滨工业大学的校训是规格严格功夫到家  
哈尔滨工业大学的校训是规格严格功夫到家  
哈尔滨工业大学的校训是规格严格功夫到家  
哈尔滨工业大学的校训是规格严格功夫到家  
哈尔滨工业大学的校训是规格严格功夫到家

### 图片版（二值 8 位灰度图）

哈	尔	滨	工	业	大	学	的	校	训	是	规	格	严	格	功	夫	到	家
L_1	L_2	L_3	L_4	L_5	L_6	L_7	L_8	L_9	L_10	L_11	L_12	L_13	L_14	L_15	L_16	L_17	L_18	L_19
哈	尔	滨	工	业	大	学	的	校	训	是	规	格	严	格	功	夫	到	家
L_20	L_21	L_22	L_23	L_24	L_25	L_26	L_27	L_28	L_29	L_30	L_31	L_32	L_33	L_34	L_35	L_36	L_37	L_38
哈	尔	滨	工	业	大	学	的	校	训	是	规	格	严	格	功	夫	到	家
L_39	L_40	L_41	L_42	L_43	L_44	L_45	L_46	L_47	L_48	L_49	L_50	L_51	L_52	L_53	L_54	L_55	L_56	L_57
哈	尔	滨	工	业	大	学	的	校	训	是	规	格	严	格	功	夫	到	家
L_58	L_59	L_60	L_61	L_62	L_63	L_64	L_65	L_66	L_67	L_68	L_69	L_70	L_71	L_72	L_73	L_74	L_75	L_76
哈	尔	滨	工	业	大	学	的	校	训	是	规	格	严	格	功	夫	到	家
L_77	L_78	L_79	L_80	L_81	L_82	L_83	L_84	L_85	L_86	L_87	L_88	L_89	L_90	L_91	L_92	L_93	L_94	L_95
哈	尔	滨	工	业	大	学	的	校	训	是	规	格	严	格	功	夫	到	家
L_96	L_97	L_98	L_99	L_100	L_101	L_102	L_103	L_104	L_105	L_106	L_107	L_108	L_109	L_110	L_111	L_112	L_113	L_114
哈	尔	滨	工	业	大	学	的	校	训	是	规	格	严	格	功	夫	到	家
L_115	L_116	L_117	L_118	L_119	L_120	L_121	L_122	L_123	L_124	L_125	L_126	L_127	L_128	L_129	L_130	L_131	L_132	L_133
哈	尔	滨	工	业	大	学	的	校	训	是	规	格	严	格	功	夫	到	家
L_134	L_135	L_136	L_137	L_138	L_139	L_140	L_141	L_142	L_143	L_144	L_145	L_146	L_147	L_148	L_149	L_150	L_151	L_152

测试集

文字版

略

图片版（二值 8 位灰度图）

哈 T_1 哈 T_20 哈 T_39 哈 T_58	尔 T_2 尔 T_21 尔 T_40 尔 T_59	滨 T_3 滨 T_22 滨 T_41 滨 T_60	工 T_4 工 T_23 工 T_42 工 T_61	业 T_5 业 T_24 业 T_43 业 T_62	大 T_6 大 T_25 大 T_44 大 T_63	学 T_7 学 T_26 学 T_45 学 T_64	的 T_8 的 T_27 的 T_46 的 T_65	校 T_9 校 T_28 校 T_47 校 T_66	训 T_10 训 T_29 训 T_48 训 T_67	是 T_11 是 T_30 是 T_49 是 T_68	规 T_12 规 T_31 规 T_50 规 T_69	格 T_13 格 T_32 格 T_51 格 T_70	严 T_14 严 T_33 严 T_52 严 T_71	格 T_15 格 T_34 格 T_53 格 T_72	功 T_16 功 T_35 功 T_54 功 T_73	夫 T_17 夫 T_36 夫 T_55 夫 T_74	到 T_18 到 T_37 到 T_56 到 T_75	家 T_19 家 T_38 家 T_57 家 T_76
哈 T_77 哈 T_96 哈 T_115 哈 T_134	尔 T_78 尔 T_97 尔 T_116 尔 T_135	滨 T_79 滨 T_98 滨 T_117 滨 T_136	工 T_80 工 T_99 工 T_118 工 T_137	业 T_81 业 T_100 业 T_119 业 T_138	大 T_82 大 T_101 大 T_120 大 T_139	学 T_83 学 T_102 学 T_121 学 T_140	的 T_84 的 T_103 的 T_122 的 T_141	校 T_85 校 T_104 校 T_123 校 T_142	训 T_86 训 T_105 训 T_124 训 T_143	是 T_87 是 T_106 是 T_125 是 T_144	规 T_88 规 T_107 规 T_126 规 T_145	格 T_89 格 T_108 格 T_127 格 T_146	严 T_90 严 T_109 严 T_128 严 T_147	格 T_91 格 T_110 格 T_129 格 T_148	功 T_92 功 T_111 功 T_130 功 T_149	夫 T_93 夫 T_112 夫 T_131 夫 T_150	到 T_94 到 T_113 到 T_132 到 T_151	家 T_95 家 T_114 家 T_133 家 T_152