

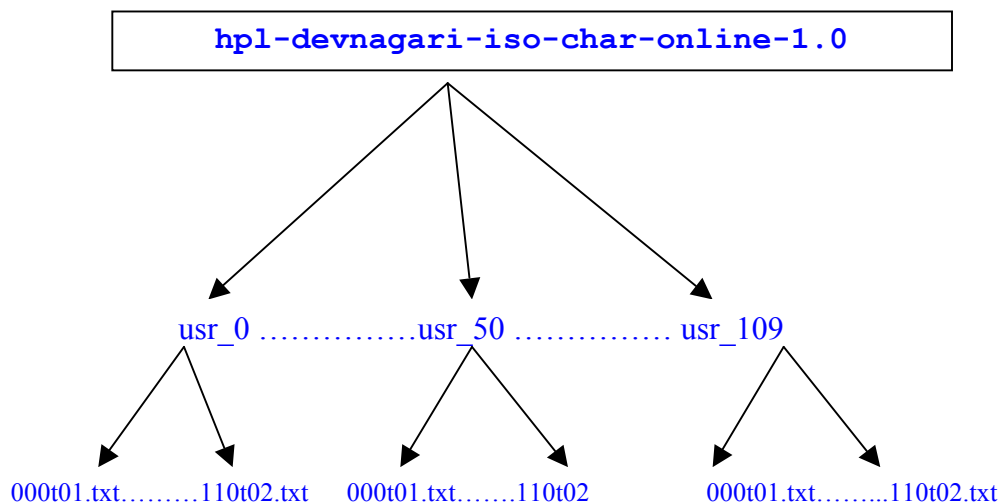
# Dataset *hpl-devnagari-iso-char-online-1.0* and *hpl-devnagari-iso-char-offline-1.0*

Updated Jun 9, 2009

The dataset *hpl-devnagari-iso-char-online-1.0* contains samples of the 111 character classes collected from different writers on ACECAD DigiMemo (A4 sized) using a *AcecadDigimemoDCT* application. The data collection forms were designed using the Adobe Designer, and were printed on A4 sheets. The writers were asked to write on the printed A4 sheets using ACECAD DigiMemo. Most writers contributed two samples per class (“trials”). Details of the directory structure, file contents and statistics regarding the number of samples per class are provided in the following sections.

The directory structure is same for both the online and offline datasets except that the name of the root directory in the offline case would be *hpl-devnagari-iso-char-offline-1.0*. The offline data contains bilevel TIFF images derived from the online data.

## 1. Directory Structure



- *hpl-devnagari-iso-char-online-1.0* directory is the root directory for the online data.
- Ink data is organized by writer into subdirectories of the form *usr\_<user-id>*, e.g., *usr\_16* corresponds to the 17<sup>th</sup> writer.
  - User-ids are in the range 0...109, but not contiguous. Samples of four writers (*usr\_9*, *usr\_58*, *usr\_72*, and *usr\_95*) were found to be very noisy and were discarded.
- Ink data is stored in files of the form *<3-digit class-id>t<trial-id>*, e.g. *008t03.txt* implies the 3<sup>rd</sup> trial of the character with class-id 008.
  - Class-id is in the range 000 ...110
  - Trial-id is in the range 01 ... 02. However, some samples written by a user might have been found noisy and hence discarded.

## 2. Ink File Contents

- In each file, ink data is represented in UNIPEN v1.0 format, as shown below.
- The channels reported for each ink point are X,Y and T. Files corresponding to some users have valid T (time) values for the first and last points of each stroke, with intermediate values set to 0. For other users, the time channel is set to 0 for all points.

```
.VERSION 1.0
.HIERARCHY CHARACTER
.COORD X Y T
.SEGMENT CHARACTER
.X_POINTS_PER_INCH 2500
.Y_POINTS_PER_INCH 2500
.POINTS_PER_SECOND 1200
.PEN_DOWN
935 523 0
935 523 0
935 523 0
935 520 0
935 517 0
935 514 0
935 511 0
.PEN_UP
```

For more information on UNIPEN format please refer to <http://unipen.nici.ru.nl/unipen.def>

## 3. Samples per class

Character Id.	Number of Samples
0	201
1	205
2	206
3	202
4	205
5	203
6	204
7	199
8	204
9	204
10	204
11	202
12	204
13	203
14	206
15	202
16	203
17	207
18	199
19	201
20	205

<b>21</b>	202
<b>22</b>	206
<b>23</b>	204
<b>24</b>	203
<b>25</b>	205
<b>26</b>	205
<b>27</b>	203
<b>28</b>	199
<b>29</b>	202
<b>30</b>	197
<b>31</b>	203
<b>32</b>	205
<b>33</b>	202
<b>34</b>	203
<b>35</b>	204
<b>36</b>	208
<b>37</b>	206
<b>38</b>	208
<b>39</b>	208
<b>40</b>	205
<b>41</b>	202
<b>42</b>	206
<b>43</b>	205
<b>44</b>	205
<b>45</b>	201
<b>46</b>	200
<b>47</b>	204
<b>48</b>	199
<b>49</b>	205
<b>50</b>	203
<b>51</b>	205
<b>52</b>	206
<b>53</b>	206
<b>54</b>	207
<b>55</b>	206
<b>56</b>	206
<b>57</b>	207
<b>58</b>	196
<b>59</b>	202
<b>60</b>	203
<b>61</b>	203
<b>62</b>	203
<b>63</b>	199
<b>64</b>	203
<b>65</b>	193
<b>66</b>	200
<b>67</b>	199
<b>68</b>	199

<b>69</b>	198
<b>70</b>	202
<b>71</b>	198
<b>72</b>	194
<b>73</b>	200
<b>74</b>	196
<b>75</b>	176
<b>76</b>	197
<b>77</b>	207
<b>78</b>	199
<b>79</b>	197
<b>80</b>	186
<b>81</b>	203
<b>82</b>	197
<b>83</b>	202
<b>84</b>	202
<b>85</b>	205
<b>86</b>	205
<b>87</b>	202
<b>88</b>	196
<b>89</b>	199
<b>90</b>	190
<b>91</b>	206
<b>92</b>	205
<b>93</b>	193
<b>94</b>	199
<b>95</b>	198
<b>96</b>	184
<b>97</b>	205
<b>98</b>	193
<b>99</b>	205
<b>100</b>	197
<b>101</b>	204
<b>102</b>	206
<b>103</b>	185
<b>104</b>	201
<b>105</b>	193
<b>106</b>	189
<b>107</b>	196
<b>108</b>	190
<b>109</b>	179
<b>110</b>	190