

1
2
3
4
5 **Single-molecule data file format (SMD): A generalized storage format for raw and**
6 **processed single-molecule data**
7

8
9
10
11 Max Greenfeld^{a,b,‡}, Jan-Willem van de Meent^{c,‡}, Dmitri S. Pavlichin^d, Hideo Mabuchi^e,
12 Chris H. Wiggins^f, Ruben L. Gonzalez Jr.^g, and Daniel Herschlag^{a,b,h}
13

14
15 Departments of ^aChemical Engineering, ^bBiochemistry, ^dPhysics and ^eApplied Physics,
16 Stanford University, Stanford, CA 94305
17

18 Departments of ^cStatistics, ^fApplied Physics and Applied Mathematics, and ^gChemistry,
19 Columbia University, New York, NY 10027
20

21 [‡]These authors contributed equally
22

23
24 ^gCorresponding author:

25 Dept. of Chemistry
26 Columbia University
27 New York, NY 10027-3126
28 212-854-1096 (telephone)
29 212-932-1289 (fax)
30 rlg2118@columbia.edu
31

32 ^hCorresponding author:

33 Dept. of Biochemistry, B400
34 Stanford University
35 Stanford, CA, 94305-5307
36 650-723-9442 (telephone)
37 650-723-6783 (fax)
38 herschla@stanford.edu
39

40
41 Additional email addresses of authors:

42 Max Greenfeld: max.greenfeld@gmail.com
43 Jan-Willem van de Meent: janwillem.vandemeent@gmail.com
44 Dmitri S. Pavlichin: dmitrip@stanford.edu
45 Hideo Mabuchi: hmabuchi@stanford.edu
46 Chris H. Wiggins: chris.wiggins@columbia.edu

47 **Comprehensive Description of SMD**

48

49 The single-molecule data format (SMD) relies on the JavaScript Object Notation
50 (JSON) format to enforce the structure. That is, every SMD format must be a valid
51 JSON file. The JSON format is a highly generalized format that is capable of storing
52 virtually any type of data. The SMD format simply imposes an order to how the JSON
53 format should be used for storing single-molecule data.

54

55 Figure S1A displays the layout of the SMD format. There are three levels of
56 hierarchy on the SMD format. The top layer of an SMD file stores an arbitrary
57 number of uniquely identified ‘experiments.’ This allows an investigator or a
58 software package using the SMD format a large amount of flexibility to combine
59 datasets with potential diverse attributes defining those datasets in one file type.
60 The second layer of the SMD format is a group that is comprised of *traces* that have a
61 large number of attributes in common. While users of the SMD format can choose to
62 use this layer of organization as they see fit, a natural way to consider this level is
63 that it would be composed of many traces arising from the same experiment. The
64 third layer is composed of parameters that are derived from, or are unique to, a
65 particular trace.

66

67 The top-level structure of the SMD format contains four fields:

68

- **desc.** This field serves to provide a simple descriptor of the data set contained herein.

69

70

- **id.** This field serves as a unique identifier for the particular set of traces that are grouped in this data structure. By default, a MD5 algorithm is used to generate a 32 digit hexadecimal number that is practically unique. This helps to ensure that when datasets generated at different times are combined, it remains easy to track the source of each dataset.

71

72

73

74

- **attr.** The attributes field stores information related to a particular group of traces. This could be information such as the day the experiment was completed, the exact experimental conditions, or any other information that relates to the data set as a whole.

75

76

77

78

- **types.** Holds type identifiers for the **index** and **values** fields. Each field of data being stored in the **values** field should be specified here. These identifiers are ‘bool’, ‘float’, ‘double’, ‘int’, ‘long’ and ‘string’.

79

80

81

82

- **data.** Holds a list of entries for each trace, which themselves contain a set of fields:

83

84

85

- **id.** Holds a trace-specific identifier. By default a MD5 hash of the **values** structure is used.

86

87

88

- **index.** This field contains a list of row labels for the **values** matrix, which typically hold the measurement times. This field should have the same length as the data in the **values** field.

89

90

91

- **values.** This field contains the actual single-molecule data. Most simply, each data type being used is stored in a field with a descriptive name (*e.g.*, channel1). While this is primarily intended to store raw single-molecule

92 data, it could equally well be used to store window-averaged data,
93 thresholded data, fits of the data or an arbitrary number of other series
94 data.
95 • **attr.** This attributes field has much the same role as the top-level
96 attributes field, but is specific to this particular trace. Within this data
97 field a user can store any additional information they are interested in
98 storing. This could be anything from a kinetic or thermodynamic
99 parameter algorithmically determined for a particular trace to an
100 observation of that particular trace that an experimentalist wants to note
101 for future reference.

102 **SMD Toolboxes for Matlab™ and Python**

103 To facilitate the adoption of the SMD format, Matlab™ and Python toolboxes were
104 created. These toolboxes allow SMD files to be generated and validated as well as
105 simple functions such as sorting and merging to be performed. The toolboxes can be
106 found at the following link: <https://smddata.github.io/>.

107 Figure S2A depicts how data fields should be populated in a Matlab™ structure to
108 enable conversion to an SMD file using the SMD Matlab™ toolbox. The fields have the
109 same names and meaning as described above. For data formatted in this manner,
110 the following functions in the SMD toolbox can be used to manipulate these
111 structures and crate SMD files. In addition to a basic version of Matlab™, this toolbox
112 requires JSONlab, a freely available toolbox for encoding and decoding JSON files in
113 Matlab™ (see: <http://iso2mesh.sourceforge.net/cgi-bin/index.cgi?jsonlab>).

114 **smd.create(data, types, varargin):** Creates a SMD structure from supplied data.

115 **smd.write_json(filename, dataset):** Saves a SMD structure as JSON (.json) or
116 compressed JSON (.json.gz).

117 **smd.read_json(filename):** Loads a SMD structure from JSON (.json) or compressed
118 JSON (.json.gz).

119 **smd.isvalid(dataset):** Checks if supplied struct is a valid SMD instance.

120 **smd.filter(dataset):** Returns a filtered dataset by matching `id` and `attr` values, or by
121 applying a custom function with boolean output to each trace.

122 **smd.merge(data1, data2, ...):** Returns a merged dataset containing all traces in
123 multiple datasets.

124 Figure S2B depicts how data fields should be populated in a Python dictionary to
125 enable conversion to an SMD file using the SMD Python toolbox. The fields have the

129 same names and meaning as described above. For data formatted in this manner the
130 following functions in the SMD toolbox can be used to manipulate these structures
131 and crate SMD files. The functionality of the Python toolbox mimics that of the
132 Matlab™ toolbox previously described.
133

134
135

```
{
  "desc": "dataset descriptor",
  "id": "dataset hash",
  "attr": {
    "data_package": "ACME Analysis",
    "types": {
      "index": "format specifier",
      "values": {"col1": "format specifier",
                 "col2": "format specifier",
                 "col3": "format specifier",...}},
    "data": [
      {
        "id": "32 character random hash tag",
        "index": [N0, N1, N2, ... NT],
        "values": [
          "col1": [X0, X1, X2, ... XT],
          "col2": [Y0, Y1, Y2, ... YT],
          "col3": [Z0, Z1, Z2, ... ZT],...],
        "attr": {
          "variable name 1": String or Number,
          "variable name 2": String or Number,...}
      },
      {
        "id": ...
        "index": ...
        "values": ...
        "attr": ...
      },
      ...
    ]
  }
}
```

136
137
138
139
140

Figure S1: Generic outline of the SMD format in JSON notation. All fields required by the SMD notation are shown.

<p>A) <code>dataset</code> : struct .desc: string dataset descriptor .id : string 32 character dataset hash tag .attr : struct Dataset level features, such as experimental conditions, description, etc. .types: struct .index: string – data type e.g. 'int' or 'double' .values.col1: string – data type e.g. 'int' or 'double' .values.col2: string – data type e.g. 'int' or 'double' .values.col3: string – data type e.g. 'int' or 'double'data : 1 x N struct .id : string 32 character values hash tag .index: T x 1 numeric Observation time range .values: struct .col1: sm data .col2: sm data .col3: sm dataattr : struct Any trace-specific features that are not time series</p>	<p>B) <code>dataset</code> : dict ['desc']: string dataset descriptor ['id'] : string 32 character dataset hash tag ['attr'] : dict Dataset level features, such as experimental conditions, description, etc. ['types']: dict ['index']: string – data type e.g. 'int' or 'double' ['values']['col1']: string – data type e.g. 'int' or 'double' ['values']['col2']: string – data type e.g. 'int' or 'double' ['values']['col3']: string – data type e.g. 'int' or 'double' ... ['data'][i] : 1 x N dict ['id'] : string 32 character values hash tag ['index']: T x 1 numeric Observation time range ['values']: dict ['col1']: sm data ['col2']: sm data ['col3']: sm data ... ['attr'] : dict Any trace-specific features that are not time series</p>
--	--

141
142
143
144
145
146
147

Figure S2: Schematic representations of the SMD format. **(A)** Representation of single-molecule data in Matlab™ prior to conversion to the SMD format with the SMD Matlab™ toolbox. **(B)** Representation of single-molecule data in a Python dictionary prior to conversion to the SMD format with the SMD Python toolbox.