

An elaborate data set on human gait and the effect of mechanical perturbations

Jason K. Moore¹, Sandra K. Hnat¹, and Antonie J. van den Bogert¹

¹Mechanical Engineering, Cleveland State University, Cleveland, Ohio, USA, 44115.
j.k.moore19@csuohio.edu, s.hnat@vikes.csuohio.edu, a.vandenbogert@csuohio.edu

ABSTRACT

Here we share a rich gait data set collected from fifteen subjects walking at three speeds on an instrumented treadmill. Each trial consists of 120 seconds of normal walking and 480 seconds of walking while being longitudinally perturbed during each stance phase with pseudo-random fluctuations in the speed of the treadmill belt. A total of approximately 1.5 hours of normal walking (> 5000 gait cycles) and 6 hours of perturbed walking (> 20,000 gait cycles) is included in the data set. We provide full body marker trajectories and ground reaction loads in addition to a presentation of processed data that includes gait events, 2D joint angles, angular rates, and joint torques along with the open source software used for the computations. The protocol is described in detail and supported with additional elaborate meta data for each trial. This data can likely be useful for validating or generating mathematical models that are capable of simulating normal periodic gait and non-periodic, perturbed gaits.

Keywords: gait, data, perturbation

INTRODUCTION

The collection of dynamical data during human walking has a long history beginning with the first motion pictures and now with modern marker based motion capture techniques and high fidelity ground reaction load measurements. Even though years of data on thousands of subjects now exist, this data is not widely disseminated, well organized, nor available with few or no restrictions. David Winter's published normative gait data, Winter (1990), is widely used in biomechanical studies, yet it comes from relatively few subjects and only a small number of gait cycles per subject. This small source has successfully inspired many other studies, such as powered prosthetic control design, Sup et al. (2008), but success in other research fields using large sets of data for discovery lead one to believe that more elaborate data sets may benefit the field of human motion studies. To enable such work, biomechanical data needs to be shared extensively, organized, and curated to enabled future analysts.

There are some notable gait data sets and databases besides Winter's authoritative set that are publicly available. The International Society of Biomechanics has maintained a web page (<http://isbweb.org/data>) since approximately 1995 that includes data sets for download and mostly unencumbered use. For example, Vaughn, et. al's data, Vaughan et al. (1992), with kinematics and force plate measurement from several subjects is available on the site. At another website, the CGA Normative Gait Database, Kirtley (2014), Chris Kirtley shares normative gait data from

27 several studies and these files have influenced other studies, for example the average gait cycles
28 from children used in van den Bogert (2003).

29 Chester et. al, Chester et al. (2007), report on a large gait database comparison where one
30 database contained kinematic data of 409 gait cycles of children from 1 to 7 years old but the data
31 does not seem to be publicly available. This is unfortunately typical. But Tirosh et. al, recognized
32 the need for a comprehensive data base for clinical gait data and created the Gaitabase, Tirosh et al.
33 (2010). This database may contain a substantial amount of data but it is encumbered by a very
34 complicated and restrictive license and sharing scheme. However, there are examples of data with
35 less restrictions. The University of Wisconsin at LaCrosse has an easily accessible normative gait
36 data set, Willson and Kernozek (2014), from 25 subjects with lower extremity marker data from
37 multiple gait cycles and force plate measurements from a single gait cycle.

38 More recent examples of biomechanists sharing their data alongside publications are: van den
39 Bogert et al. (2013) which includes full body joint kinematics and kinetics from eleven subjects
40 walking for a small number of gait cycles and Wang and Srinivasan (2014) who includes a larger set
41 of data from ten subjects walking for five minutes each at three different speeds but only a small
42 set of lower extremity markers are present. The second is notable because it publishes the data in
43 Dryad, a modern citable data repository.

44 The publicly available gait data is small compared to the number of gait studies that have been
45 performed over the years. The data that is available generally suffers from limitations such as few
46 subjects, few gait cycles, few markers, highly clinical, no raw data, limited force plate measurements,
47 lack of meta data, non-standard formats, and restrictive licensing. To help with this situation we are
48 making the data we collected for our research purposes publicly available and free of the previously
49 mentioned deficiencies. Not only do we provide a larger set of normative gait data that has been
50 previously available, we also include an even larger set of data in which the subject is being perturbed,
51 something that does not currently exist. We believe both of these sets of data can serve a variety of
52 use cases and hope that we can save time and effort for future researchers by sharing it.

53 Our use case for the data is centered around the need of bio-inspired control systems for emerging
54 powered prosthetics and orthotics. Ideally, a powered prosthetic would behave in such a way that the
55 user would feel like their limb was never disabled. There are a variety of approaches to developing
56 bio-inspired control systems, some of which aim to mimic the reactions and motion of an able-
57 bodied person. A modern gait lab is able to collect a variety of kinematic, kinetic, and physiological
58 data from humans during gait. This data can potentially be used to drive the design of the human-
59 mimicking controller. With a rich enough data set, one may be able to identify control mechanisms
60 used during a human's natural gait and recovery from perturbations. We have collected data that is
61 richer than previous gait data sets and may be rich enough for control identification. The data can
62 also be used for verification purposes for controllers that have been designed in other manners.

63 With all of this in mind, we collected over seven and half hours of gait data from fifteen able
64 bodied subjects which amounts to over 25,000 gait cycles. The subjects walked at three different
65 speeds on an instrumented treadmill while we collected full body marker locations and ground
66 reaction loads from a pair of force plates. The protocol for the majority of the trials included two
67 minutes of normal walking and eight minutes of walking under the influence of pseudo-random belt
68 speed fluctuations. The data has been organized complete with rich meta data and made available in
69 the most unrestrictive form for other research uses following modern best practices in data sharing,
70 White et al. (2013).

71 Furthermore, we include a small Apache licensed open source software library for basic gait

Table 1. Information about the 15 participants. The final three columns give the trial numbers associated with each nominal treadmill speed. The measured mass is computed from the mean total vertical ground reaction force just after the calibration pose event, if possible. Additional trials found in the data set with a subject identification number 0 are trials with no subject, i.e. unloaded trials that can be used for inertial compensation purposes, and are not shown in the table. Generated by `src/subject_table.py`.

Id	Gender	Age [yr]	Height [m]	Measured Mass [kg]	Self-reported mass [kg]	0.8 m/s	1.2 m/s	1.6 m/s
1	male	25	1.87	NA	101	NA	6, 7, 8	NA
3	female	32	1.62	54 ± 2	60	46	47	48
4	male	30	1.76	NA	74	12, 15	13	14
5	male	23	1.73	71.2 ± 0.9	65	32	31	33
6	male	26	1.77	86.8 ± 0.6	80	40	41	42
7	female	29	1.72	64.5 ± 0.8	63	16	17	18
8	male	20	1.57	74.9 ± 0.9	70	19	20	21
9	male	20	1.69	67 ± 2	64	25	26	27
10	male	19	1.77	92 ± 2	91	61	62	63
11	male	22	1.85	NA	80	9	10	11
12	male	22	1.85	74.2 ± 0.5	81	49	50	51
13	female	21	1.70	58 ± 2	64	55	56	57
15	male	22	1.83	80.5 ± 0.8	79	67	68	69
16	female	28	1.69	56.2 ± 0.6	52	76	77	78
17	male	23	1.86	88.3 ± 0.8	87	73	74	75

72 analysis and demonstrate its use in the paper. The combination of the open data and open software
 73 allow the results presented within to be computationally reproducible and instructions are included
 74 in the associated repository for doing so.

75 METHODS

76 Participants

77 Fifteen able bodied subjects including four females and eleven males with an average age of 24 ± 4
 78 years, height of 1.75 ± 0.09 m, mass of 74 ± 13 kg participated in the study. The study was approved
 79 by the Institutional Review Board of Cleveland State University (# 29904-VAN-HS) and written
 80 informed consent was obtained from all participants. The data has been anonymized with respect
 81 to the participants' identities and a unique identification number was assigned to each subject. A
 82 selection of the meta data collected for each subject is shown in Table 1.

83 Equipment

84 The data were collected in the Laboratory for Human Motion and Control at Cleveland State
 85 University, using the following equipment:

- 86 • A R-Mill treadmill which has dual 6 degree of freedom force plates, independent belts for
 87 each foot, and lateral/pitch motion capabilities (Forcelink, Culemborg, Netherlands).

- 88 • A 10 Osprey camera motion capture system paired with the Cortex 3.1.1.1290 software
89 (Motion Analysis, Santa Rosa, CA, USA).
- 90 • USB-6255 data acquisition unit (National Instruments, Austin, Texas, USA).
- 91 • Four ADXL330 Triple Axis Accelerometer Breakout boards attached to the treadmill (Spark-
92 fun, Niwot, Colorado, USA).
- 93 • D-Flow software (versions 3.16.1 to 3.16.2) and visual display system, (Motek Medical,
94 Amsterdam, Netherlands).

95 The Cortex software delivers high accuracy 3D marker trajectories from the cameras along with
96 data from force plates and analog sensors (EMG/Accelerometer) through a National Instruments
97 USB-6255 data acquisition unit. D-Flow is required to collect data from any digital sensors and to
98 control the treadmill's motion (lateral, pitch, and belts). D-Flow can process the data in real time
99 and/or export data to file.

100 Our motion capture system's coordinate system is such that the X coordinate points to the right,
101 the Y coordinate points upwards, and the Z coordinate follows from the right-hand-rule, i.e. points
102 backwards with respect to the walking direction. The camera's coordinate system is aligned to an
103 origin point on treadmill's surface during camera calibration. The same point is used as the origin of
104 the ground reaction force measuring system. Figure 1 shows the layout of the equipment.

105 Early on, we discovered that the factory setup of the R-Link treadmill had a vibration mode as
106 low as 5Hz that is detectable in the force measurements, likely due to the flexible undercarriage
107 and pitch motion mechanism. Trials 6-8 are affected by this vibration mode. During trials 9-15
108 the treadmill was stabilized with wooden blocks. During, the remaining trials the treadmill was
109 stabilized with metal supports. See the Data Limitations Section for more details.

110 The acceleration of the treadmill was measured during each trial by four ADXL330 accelerom-
111 eters placed at the four corners of the machine. These accelerometers were intended to provide
112 information for inertial compensation purposes when the treadmill moved laterally, but are extrane-
113 ous for trials greater than 8 due to the treadmill being stabilized.

114 Protocol

115 The experimental protocol consisted of both static measurements and walking on the treadmill for
116 10 minutes under unperturbed and perturbed conditions. Before a set of trials on the same day the
117 following happened:

- 118 • Calibration of the motion capture system using the manufacturer's recommended procedure.
- 119 • Subject changes into athletic shoes, shorts, sports bra, baseball cap, and rock climbing harness.
- 120 • All 47 markers are applied directly to the skin except for the heel, toe, and head markers,
121 which were placed on the respective article of clothing.¹
- 122 • Subjects self-reported age, gender, and mass.
- 123 • Height was measured by the experimentalist.

¹The sacrum and rear pelvic markers may have been placed on the shorts for a small number of the subjects



Figure 1. The treadmill with coordinate system, cameras (circled in orange), projection screen, and safety rope. The direction of travel is in the $-z$ direction.

- Four reference photographs (front, back, right, left) were taken of subject's marker locations.

After obtaining informed consent and a briefing by the experimentalist on the trial protocol, the subject followed the verbal instructions of the experimentalist and the on-screen instructions from the video display. The protocol for a single trial was as follows:

1. Subject stepped onto the treadmill and markers were identified with Cortex.
2. The safety rope was attached loosely to the rock climbing harness such that no undue forces were acting on the subject during walking, but that the harness would prevent a full fall.
3. The subject started by stepping on sides of treadmill so that feet did not touch the force plates and the force plate signals are zeroed. This corresponds to the "Force Plate Zeroing" event.
4. Once notified by the video display, the subject stood in the initialization pose: standing straight up, looking forward, arms out by their sides (45 degrees) and the event, "Calibration Pose", was manually recorded by the operator.
5. A countdown to the first normal walking phase was displayed. At the end of the countdown the event "First Normal Walking" was recorded and the treadmill ramped up to the specified speed and the subject was instructed to walk normally, to focus on the "endless" road on the display, and not to look at their feet.
6. After 1 minute of normal walking, the longitudinal perturbation phase begun and was recorded as "Longitudinal Perturbation".
7. After 8 minutes of walking under the influence of the perturbations, the second normal walking phase begun and was recorded as "Second Normal Walking".

- 144 8. After 1 minute of normal walking, a countdown was shown on the display and the treadmill
145 decelerated to a stop.
- 146 9. The subject was instructed to step off of the force plates for 10 seconds and the “Unloaded
147 End” event was recorded.
- 148 10. The subject could then take a rest break before each additional trial.

149 Trials 6-8 included a calibration pose at the start of the trial but the event was not explicitly
150 recorded. In those trials, the “TreadmillPerturbation” event marks the beginning of longitudinal
151 perturbations and the “Both” event marks the beginning of combined longitudinal and lateral
152 perturbations. The force plate zeroing at the end was also not explicitly recorded.

153 Perturbation Signals

154 As previously described, the protocol included a phase of normal walking, followed by longitudinal
155 belt speed perturbations, and ended with a second segment of normal walking. Three pseudo-
156 random belt speed control signals, with mean velocities of 0.8 m s^{-1} , 1.2 m s^{-1} and 1.6 m s^{-1} , were
157 pre-generated with MATLAB and Simulink (Mathworks, Natick, Massachusetts, USA). The same
158 control signal was used for all trials at that given speed.

159 To create the signals, we started by generating random 100 Hz acceleration signals using the
160 Simulink discrete-time Gaussian white noise block followed by a saturation block set at the maximum
161 belt acceleration of 15 m s^{-2} . The signal was then integrated to obtain belt speed and high-pass
162 filtered with a second-order Butterworth filter to eliminate drift. One of the three mean speeds
163 were then added to the signal and limited between 0 m s^{-1} to 3.6 m s^{-1} . The cutoff frequencies
164 of the high-pass filter, as well as the variance in the acceleration signal, were manually adjusted
165 until acceptable standard deviations for each mean speed were obtained: 0.06 m s^{-1} , 0.12 m s^{-1} and
166 0.21 m s^{-1} for the three speeds, respectively. These ensured that the test subjects were sufficiently
167 perturbed at each speed, while remaining within the limits of our equipment and testing protocol.

168 To ensure that the treadmill belts could accelerate to the desired values, the high performance
169 mode in the D-Flow software was enabled. This had the side effect of enabling too rapid of
170 accelerations when the belt speed changed to or from zero speed. To eliminate this, a suitable
171 ramped acceleration and deceleration were generated for the speed transitions.

172 The MATLAB script and Simulink model produce a comma-delimited text file of six signals:
173 time stamp, slow, normal, and fast walking perturbation signals, and slow and fast running signals.²
174 The measured speed of the treadmill belts are compared to the control input signals in Figure 2 to
175 show the effect of the treadmill and controller dynamics. The system introduces a delay and seems
176 to act as a low pass filter. The standard deviations of the outputs do not significantly differ from the
177 desired values: 0.05 m s^{-1} , 0.12 m s^{-1} and 0.2 m s^{-1} for the three speeds, respectively.

178 To show the effects of the treadmill dynamics and give an idea of the frequency content of
179 the actual perturbations, the input and output for each speed were transformed into the frequency
180 domain using the Fast Fourier algorithm, and the results are shown in Figure 3. This shows that for
181 the 1.2 m s^{-1} walking speed, the amplitude of the output is significantly lower than the amplitude of
182 the input signal at lower frequencies. Additionally, the amplitude of output signal in the 0.8 m s^{-1}
183 walking speed begins to attenuate around 2 Hz, which is a noticeably lower frequency than the other
184 walking speeds.

²The running signals were not used in the experiments presented in this paper.

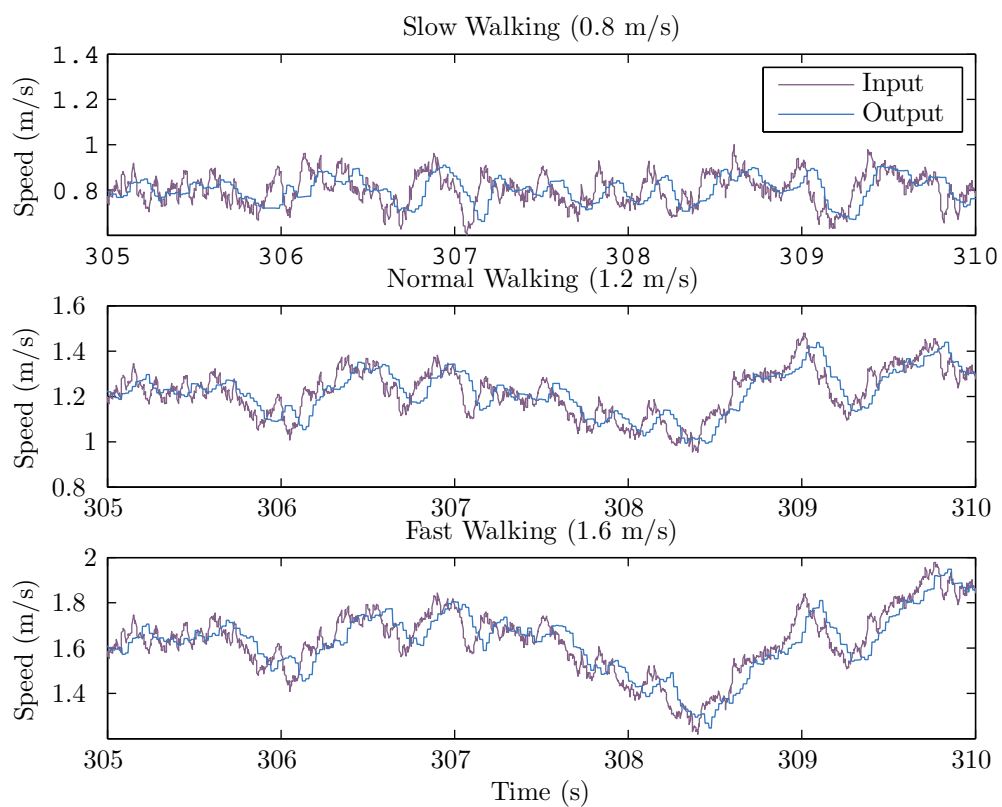


Figure 2. Treadmill belt speed input signals (purple) and recorded output speeds (blue) for average belt speeds of 0.8 m s^{-1} , 1.2 m s^{-1} and 1.6 m s^{-1} , respectively.

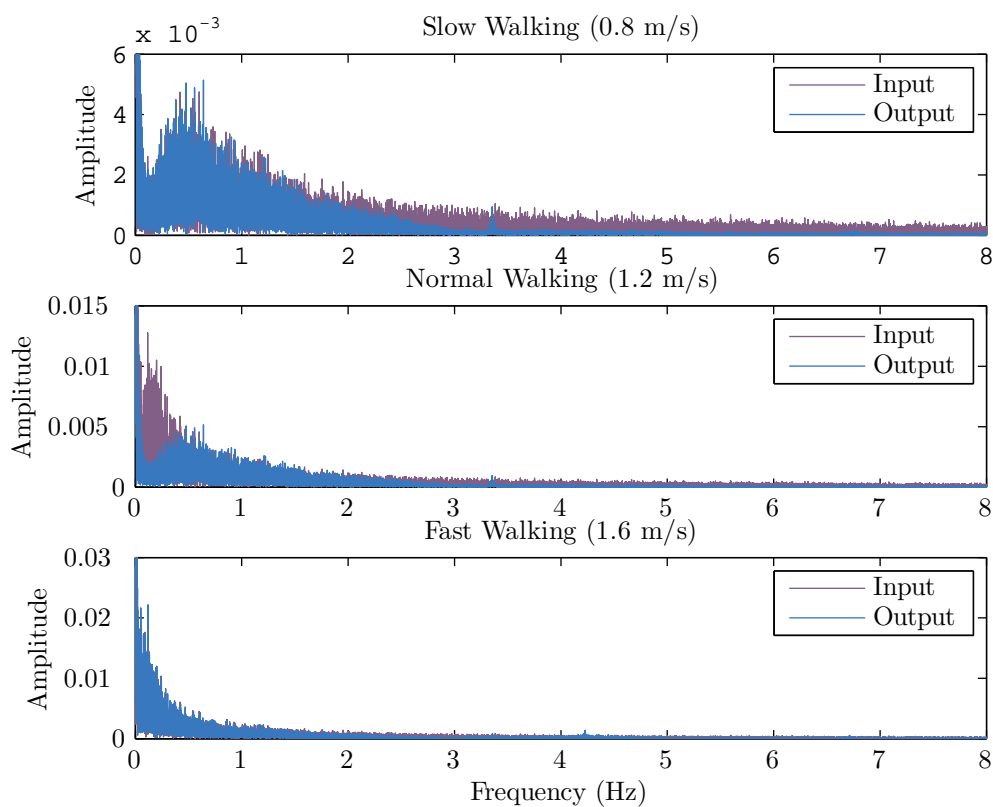


Figure 3. Frequency spectrum of the treadmill belt velocity input signal (purple) and the recorded output velocity (blue) for average belt speeds of 0.8, 1.2, and 1.6 m/s, respectively

185 RESULTS

186 Raw Data

187 The raw data consists of a set of ASCII tab delimited text files output from both the “mocap” and
188 “record” modules in D-Flow in addition to a manually generated YAML file that contains all of the
189 necessary meta data for the given trial. These three files are stored in a hierarchy of directories with
190 one trial per directory. The directories are named in the following fashion T001/ where T stands
191 for “trial” and the following three digits are provide a unique trial identification number.

192 *mocap-xxx.txt*

193 The output from the D-Flow mocap module is stored in a tab delimited file named `mocap-xxx.txt`
194 where `xxx` represents the trial id number. The file is tab delimited and contains a number of time
195 series. The numerical values of the time series are provided in decimal fixed point notation with
196 6 decimals of precision, e.g. `123456.123456`, regardless of the units. The first line of the file
197 holds the header. The header includes time stamp column, frame number column, marker position
198 columns, force plate force/moment columns, force plate center of pressure columns, other analog
199 columns, and potentially results from the real time Human Body Model van den Bogert et al. (2013)
200 which is included with D-Flow. The columns are further described below:

201 **TimeStamp** The monotonically increasing computer clock time when D-Flow receives a frame
202 from Cortex. These are recorded at approximately at 100 Hz and given in seconds.

203 **FrameNumber** Monotonically increasing positive integers that correspond to each frame received
204 from Cortex.

205 **Marker Coordinates** Any column that ends in `.PosX`, `.PosY`, or `.PosZ` are marker coordinates
206 expressed in Cortex’s Cartesian reference frame. The prefixes match the marker labels given
207 in Table 2. These values are in meters.

208 **Ground Reaction Loads** There are three ground reaction forces and three ground reaction moments
209 recorded by each of the two force plates in Newtons and Newton-Meters, respectively. The
210 prefix for these columns is either `FP1` or `FP2` and represents either force plate 1 (left) or
211 2 (right). The suffixes are either `.For[XYZ]`, `.Mom[XYZ]` for the forces and moments,
212 respectively. The force plate voltages are sampled at a much higher frequency than the
213 cameras, but delivered at the Cortex camera sample rate, 100 Hz through the D-Flow mocap
214 module. A force/moment calibration matrix stored in Cortex converts the voltages to forces
215 and moments before sending it to D-Flow. Cortex also computes the center of pressure from
216 the forces, moments, and force plate dimensions. These have the same prefixes for the plate
217 number, have the suffix `.Cop[XYZ]`, and are given in meters.

218 **Analog Channels** Several analog signals are recorded under column headers `Channel[1-99].Anlg.`
219 These correspond to analog signals sampled by Cortex and correspond to the 96 analog chan-
220 nels in the National Instruments USB-6255. The first twelve are the voltages from the force
221 plate load cells. We also record the acceleration of 4 points on the treadmill base in analog
222 channels 61-72 that were in place in case inertial compensation for the lateral treadmill
223 movement was required.

224 **record-xxx.txt**

225 The record module also outputs a tab delimited ASCII text file with numerical values at six decimal
226 digits. It includes a `Time` column which records the D-Flow system time in seconds. This time
227 corresponds to the time recorded in the `TimeStamp` column in mocap module tsv file which is
228 necessary for time synchronization. There are two additional columns `RightBeltSpeed` and
229 `LeftBeltSpeed` which provide the independent belt speeds measured in meters per second by a
230 factory installed encoder in the treadmill.

231 Additionally, the record module is capable of recording the time at which various preprogrammed
232 events occur, as detected or set by D-Flow. It does this by inserting commented (`#`) lines in between
233 the rows when the event occurred. The record files have several events that delineate the different
234 phases of the protocol:

235 **A: Force Plate Zeroing** Marks the time at the beginning of the trial at which there is no load on
236 the force plates and when the force plate voltages were zeroed.

237 **B: Calibration Pose** Marks the time at which the person is in the calibration pose.

238 **C: First Normal Walking** Marks the time when the treadmill begins Phase 1: constant belt speed.

239 **D: Longitudinal Perturbation** Marks the time when the treadmill begins Phase 2: longitudinal
240 perturbations in the belt speed.

241 **E: Second Normal Walking** Marks the time when phase 3 starts: constant belt speed.

242 **F: Unloaded End** Marks the time at which there is no load on the force plates and the belts are
243 stationary.

244 **meta-xxx.yml**

245 Each trial directory contains a meta data file in the YAML format named in the following style
246 `meta-xxx.yml` where `xxx` is the three digit trial identification number. There are three main
247 headings in the file: `study`, `subject`, and `trial`. An example meta data file is shown in Listing
248 1.

249 The `study` section contains identifying information for the overall study, an identification
250 number, name, and description. This is the same for all meta data files in the study. Details are given
251 below:

252 **id** An integer specifying a unique identification number of the study.

253 **name** A string giving the name of the study.

254 **description** A string with a basic description of the study.

255 The `subject` section provides key value pairs of information about the subject in that trial.
256 Each subject has a unique identification number along with basic anthropomorphic data. The
257 following details the possible meta data for the subject:

258 **age** An integer age in years of the subject at the time of the trial.

259 **ankle-width-left** A float specifying the width of the subjects left ankle.

260 **ankle-width-right** A float specifying the width of the subjects right ankle.

261 **ankle-width-units** A string giving the units of measurement of the ankle widths.

262 **id** An unique identification integer for the subject.

263 **gender** A string specifying the gender of the subject.

264 **height** A float specifying the measured height of the subject (with shoes and hat on) at the time of
265 the trial.

266 **height-units** A string giving the units of the height measurement.

267 **knee-width-left** A float specifying the width of the subjects left knee.

268 **knee-width-right** A float specifying the width of the subjects right knee.

269 **knee-width-units** A string giving the units of measurement of the knee widths.

270 **mass** A float specifying the self-reported mass of the subject.

271 **mass-units** A string specifying the units of the mass measurement.

272 The `trial` section contains the information about the particular trial. Each trial has a unique
273 identification number along with a variety of other information, detailed below:

274 **analog-channel-map** A mapping of the strings D-Flow assigns to signals emitted from the analog
275 channels of the NI USB-6255 to names the user desires.

276 **cortex-version** The version of Cortex used to record the trial.

277 **datetime** A date formatted string giving the date of the trial in the `YYYY-MM-DD` format.

278 **dflow-version** The version of D-Flow used to record the trial.

279 **events** A key value map which prescribes names to the alphabetic events recorded in the record file.

280 **files** A key value mapping of files associated with this trial where the key is the D-Flow file type
281 and the value is the path to the file relative to the meta file. The compensation file corresponds
282 to an unloaded trial collected on the same day that could be used for inertial compensation
283 purposes, if needed.

284 **hardware-settings** There are tons of settings for the hardware in both D-Flow, Cortex, and the
285 other software in the system. This contains any non-default settings.

286 **high-performance** A boolean value indicating whether the D-Flow high performance setting
287 was on (True) or off (False).

288 **id** An unique three digit integer identifier for the trial. All of the file names and directories associated
289 with this trial include this number.

290 **marker-map** A key value map which maps marker names in the mocap file to the user's desired
291 names for the markers.

292 **marker-set** Indicates the HBM van den Bogert et al. (2013) marker set used during the trial, either
293 full, lower, or NA.

294 **nominal-speed** A float representing the nominal desired treadmill speed during the trial.

295 **nominal-speed-units** A string providing the units of the nominal speed.

296 **notes** Any notes about the trial.

297 **pitch** A boolean that indicates if the treadmill pitch degree of freedom was actuated during the trial.

298 **stationary-platform** A boolean that indicates whether the treadmill sway or pitch motion was
299 actuated during the trial. If this flag is false, the measured ground reaction loads must be
300 compensated for the inertial affects and be expressed in the motion capture reference frame.

301 **subject-id** An integer corresponding to the subject in the trial.

302 **sway** A boolean that indicates if the treadmill lateral degree of freedom was actuated during the
303 trial.

304 **Markers**

305 We make use of the full body 47 marker set described in van den Bogert et al. (2013) and presented
306 in detail in Table 2. As with all camera based motion capture systems, the markers sometimes go
307 missing in the recording. When a marker goes missing, if the data was recorded in a D-Flow version
308 less than 3.16.2rc4 [3], D-Flow continues to record the last non-missing value in all three axes
309 until the marker is visible again. In D-Flow versions greater than or equal to 3.16.2rc4, the missing
310 markers are indicated in the TSV file as either 0.000000 or -0.000000, which is the same as
311 has been in the HBM columns in all versions of D-Flow. The D-Flow version must be provided in
312 the meta data yml file for each trial to be able to distinguish this detail.

313 **Processed Data**

314 We developed a toolkit for data processing, GaitAnalysisToolKit v0.1.2, Moore et al. (2014b), for
315 common gait computations and provide an example processed trial to present the nature of the data.
316 The tool was developed in Python, is dependent on the SciPy Stack and Octave, and provides two
317 main classes: one to do basic gait data cleaning from D-Flow's output files, `DFlowData`, and a
318 second to compute common gait variables of interest, `GaitData`.

319 The `DFlowData` class collects and stores all the raw data presented in the previous section and
320 applies several "cleaning" operations to transform the data into a usable form. The cleaning process
321 follows these steps:

- 322 1. Load the meta data file into a Python dictionary.
- 323 2. Load the D-Flow mocap module TSV file into Pandas `DataFrame`.
- 324 3. Relabel the column headers to more meaningful names if this is specified in the meta data.

Table 2. Descriptions of the 47 markers used in this study. The “Set” column indicates whether the marker exists in the lower and/or full body marker set. The label column matches the column headers in the mocap-xxx.txt files and/or the marker map in the meta-xxx.yml file.

Set	#	Label	Name	Description
F	1	LHEAD	Left head	Just above the ear, in the middle.
F	2	THEAD	Top head	On top of the head, in line with the LHEAD and RHEAD.
F	3	RHEAD	Right head	Just above the ear, in the middle.
F	4	FHEAD	Forehead	Between line LHEAD/RHEAD and THEAD a bit right from center.
L/F	5	C7	C7	On the 7th cervical vertebrae.
L/F	6	T10	T10	On the 10th thoracic vertebrae.
L/F	7	SACR	Sacrum bone	On the sacral bone.
L/F	8	NAVE	Navel	On the navel.
L/F	9	XYPH	Xiphoid process	Xiphoid process of the sternum.
F	10	STRN	Sternum	On the jugular notch of the sternum.
F	11	BBAC	Scapula	On the inferior angle fo the right scapular.
F	12	LSHO	Left shoulder	Left acromion.
F	13	LDELT	Left deltoid muscle	Apex of the deltoid muscle.
F	14	LLEE	Left lateral elbow	Left lateral epicondyle of the elbow. Upper one in the T-Pose.
F	15	LMEE	Left medial elbow	Left medial epicondyle of the elbow. Lower on in the T-Pose.
F	16	LFM	Left forearm	On 2/3 on the line between the LLEE and LMW.
F	17	LMW	Left medial wrist	On styloid process radius, thumb side.
F	18	LLW	Left lateral wrist	On styloid process ulna, pinky side.
F	19	LFIN	Left fingers	Center of the hand. Caput metatarsal 3.
F	20	RSHO	Right shoulder	Right acromion.
F	21	RDELT	Right deltoid muscle	Apex of deltoid muscle.
F	22	RLEE	Right lateral elbow	Right lateral epicondyle of the elbow. Lower one in the T-pose.
F	23	RMEE	Right medial elbow	Right medial epicondyle of the elbow. Lower one in the T-pose.
F	24	RFRM	Right forearm	On 1/3 on the line between the RLEE and RMW.
F	25	RMW	Right medial wrist	On styloid process radius, thumb side.
F	26	RLW	Right lateral wrist	On styloid process ulna, pinky side.
F	27	RFIN	Right fingers	Center of the hand. Caput metatarsal 3.
L/F	28	LASIS	Pelvic bone left front	Left anterior superior iliac spine.
L/F	29	RASIS	Pelvic bone right front	Right anterior superior iliac spine.
L/F	30	LPSIS	Pelvic bone left back	Left posterior superior iliac spine.
L/F	31	RPSIS	Pelvic bone right back	Right posterior superior iliac spine.
L/F	32	LGTRO	Left greater trochanter of the femur	On the cetner of the left greater trochanter.
L/F	33	FLTHI	Left thigh	On 1/3 on the line between the LFTRO and LLEK.
L/F	34	LLEK	Left lateral epicondyle of the knee	On the lateral side of the joint axis.
L/F	35	LATI	Left anterior of the tibia	On 2/3 on the line between the LLEK and LLM.
L/F	36	LLM	Left lateral malleolus of the ankle	The center of the heel at the same height as the toe.
L/F	37	LHEE	Left heel	Center of the heel at the same height as the toe.
L/F	38	LTOE	Left toe	Tip of big toe.
L/F	39	LMT5	Left 5th metatarsal	Caput of the 5th metatarsal bone, on joint line midfoot/toes.
L/F	40	RGTRO	Right greater trochanter of the femur	On the cetner of the right greater trochanter.
L/F	41	FRTHI	Right thigh	On 2/3 on the line between the RFTRO and RLEK.
L/F	42	RLEK	Right lateral epicondyle of the knee	On the lateral side of the joint axis.
L/F	43	RATI	Right anterior of the tibia	On 1/3 on the line between the RLEK and RLM.
L/F	44	RLM	Right lateral malleolus of the ankle	The center of the heel at the same height as the toe.
L/F	45	RHEE	Right heel	Center of the heel at the same height as the toe.
L/F	46	RTOE	Right toe	Tip of big toe.
L/F	47	RMT5	Right 5th metatarsal	Caput of the 5th metatarsal bone, on joint line midfoot/toes.

- PeerJ PrePrints
- 325 4. Optionally identify the missing values in the mocap marker data and replace them with
326 `numpy.nan`.
 - 327 5. Optionally interpolate the missing marker values and replaces them with interpolated estimates
328 using a variety of interpolation methods.
 - 329 6. Load the D-Flow record module TSV file into a Pandas `DataFrame`.
 - 330 7. Extract the events and create a dictionary mapping the event names in the meta data to the
331 events detected in the record module file.
 - 332 8. Internally compensate the ground reaction loads based on whether the meta data indicates
333 there was treadmill motion.
 - 334 9. Merge the data from the mocap module and record module into one data frame at the maximum
335 common constant sample rate.

336 Once the data is cleaned there are two methods that allow you to extract the cleaned data: either
337 extract sections of the data bounded by the events recorded in the `record-xxx.txt` file or save
338 the cleaned data to disk. These operations are available as a command line application and as an
339 application programming interface (API) in Python. An example of the `DFlowData` API in use is
340 provided in Listing 2.

341 The `GaitData` class is then used to compute things such as gait events (toe off and heel strike
342 times), basic 2D kinematics and inverse dynamics, and to store the data into a Pandas `Panel` with
343 each gait cycle on the item axis at a specified sample rate. This object can also be serialized to disk
344 in HDF5 format. An example of using the Python API is shown in Listing 3.

345 A similar work flow was used to produce Figure 4 which compares the mean and standard
346 deviation of sagittal plane joint angles and torques from the perturbed gait cycles and the unperturbed
347 gait cycles computed from trial 20. This gives an idea of the more highly variable dynamics required
348 to walk while being longitudinally perturbed.

349 For more insight into the difference in the unperturbed and perturbed data, Figure 5 compares
350 the distribution of a few gait cycle statistics. One can see that the perturbed strides have a much
351 larger variation in frequency and length.

352 Data Limitations

353 The data is provided in good faith with great attention to detail but as with all data there are anomalies
354 that may affect the use and interpretation of results emanating from the data. The following list gives
355 various notes and warnings about the data that should be taken into account when making use of it.

- 356 • Be sure to read the notes in each meta data file for details about possible anomalies in that
357 particular trial. Things such as marker dropout, ghost markers, and marker movement are the
358 more prominent notes. Details about variations in the equipment on the day of the trial are
359 also mentioned.
- 360 • The subject identification number 0 stands for "no subject" and was used whenever data was
361 collected from the system with no subject on the treadmill, for example during the trials that
362 were intended to be used for inertial compensation purposes. These trials play through the
363 exact protocol as those with a human subject and the matching trials are indicated in the meta

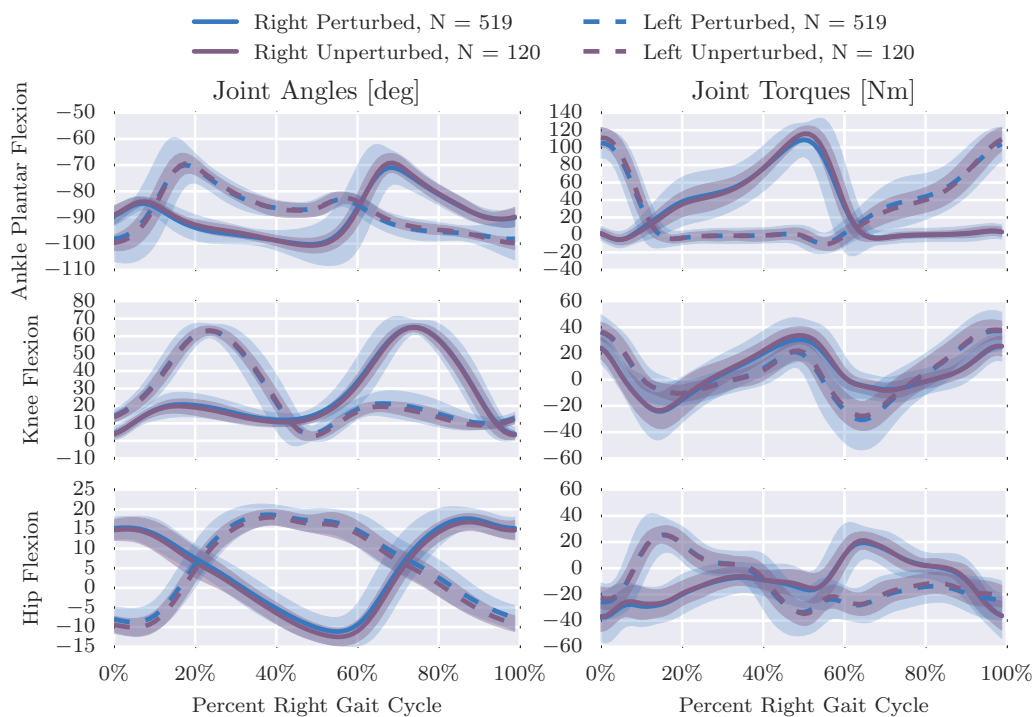


Figure 4. Mean (right: solid, left: dashed) and 3σ (shaded) joint angles and torques from both unperturbed (purple) and perturbed (blue) gait cycles from trial 20. Produced by `src/unperturbed_perturbed_comparison.py`.

data. Matching unloaded trials were recorded on the same day as the loaded trials and is noted in the `trial:files:compensation` section of the meta data file.

- Trials 1 and 2 were not recorded as part of this study. Those trial identification numbers were reserved for early data exploration from data collected in other studies.
- Trials 37, 38, and 39 do not exist. The numbers were accidentally skipped.
- Trials 9, 10, and 11 used a slightly different event definition where the calibration poses were not explicitly tagged by an event, yet the protocol was the identical to the following trials. The calibration pose will have to be determined manually.
- Trials 6-15 have force measurements are affected by the treadmill vibration mode mentioned in the equipment section and the forces should be not be used. We include the trials because both the kinematic data is valid and trials 6-8 include lateral perturbations in addition to the longitudinal.
- During trials 9-15 we used wooden blocks to fix the treadmill to the concrete floor to eliminate the treadmill's low vibration mode (~ 5 Hz). But these blocks seem to have corrupted the force plate measurements by imposing frictional stresses on the system. The force plate measurements should not be used from these trials, but the marker data is fine.
- Trials 6-8 use an early experimental protocol which divided the perturbation sections into three sections: longitudinal perturbations, lateral perturbations, and a combination of each.

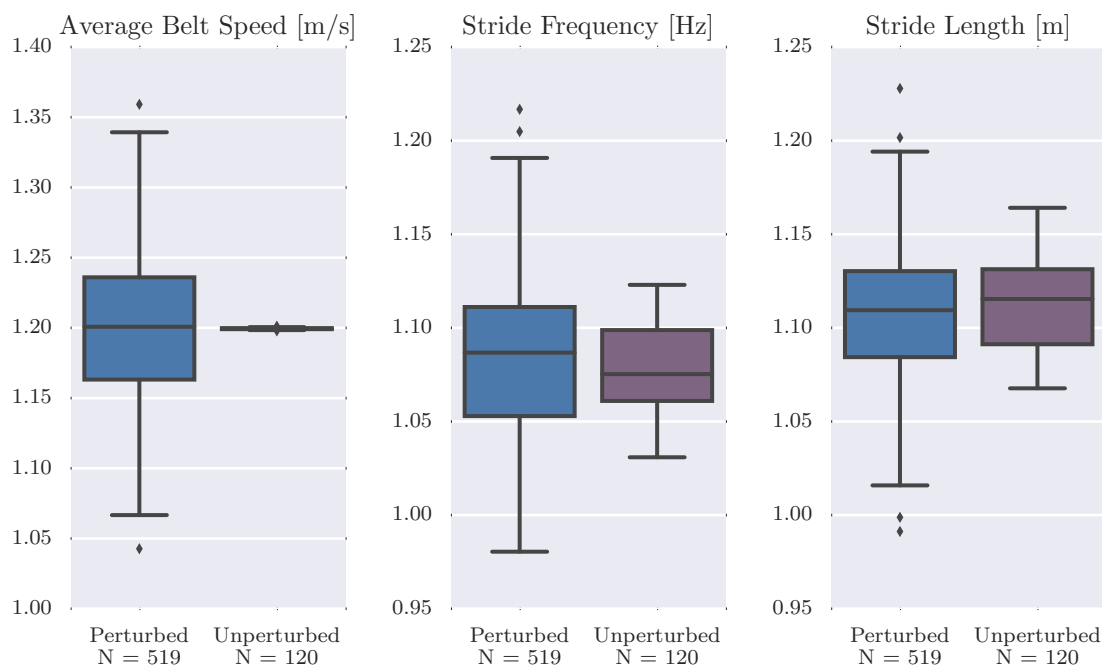


Figure 5. Box plots of the average belt speed, stride frequency, and stride length which compare gait cycles for the unperturbed (purple) and perturbed (blue). Produced by `src/unperturbed_perturbed_comparison.py`.

382 We then learned the treadmill had a low vibrational mode which significantly affects the force
 383 plate measurements, requiring us to eliminate the lateral perturbation motions. The force
 384 measurements during these trials are corrupted by this vibrational mode and should be used
 385 with caution or not at all.

- 386 • We did not record unloaded compensation trials for trials 9-15. Regardless, they would likely
 387 be useless due to the corruption from the wooden blocks.
- 388 • Trials 6-8 use a only the lower body marker set. The remaining trials are full body.
- 389 • The ankle joint torques computed from subject 9's data in trials 25-27 are abnormal and should
 390 be used with caution or not at all. We were not able to locate the source of the error, but it is
 391 likely related to the force calibration.

392 CONCLUSION

393 We have presented a rich and elaborate data set of motion and ground reaction loads from human
 394 subjects during both normal walking and when recovering from longitudinal perturbations. The raw
 395 data is provided for reuse with complete meta data. In addition to the data, we provide software that
 396 can process the data for both cleaning purposes and to produce typical sagittal plane gait variables
 397 of interest. Among other uses, we believe the dataset is ideally suited for control identification
 398 purposes. Many researchers are working on mathematical models for control in gait and this dataset
 399 provides both a way to validate these models and a source for generating them.

400 DATA AVAILABILITY

401 The data set, Moore et al. (2014a), is available via the Zenodo data repository. Two approximately
402 1.2GB gzipped tar balls contain the data and a README file with a short description of the contents.
403 The data is released under the Creative Commons CC0 license (<http://creativecommons.org/about/cc0>)
404 following best practices for sharing scientific data.

405 SOFTWARE AVAILABILITY

406 The tables and figures in the paper can be reproduced from the source repository shared on Github:
407 <https://github.com/csu-hmc/perturbed-data-paper>. Along with the source code in the repository, the
408 computations depend on version 0.1.2 of the GaitAnalysisToolKit, Moore et al. (2014b), which can
409 be downloaded from Zenodo or the Python Package Index (<http://pypi.python.org>).

410 AUTHOR CONTRIBUTIONS

411 A.v.d.B. conceived of the experiments and protocol. J.K.M and S.K.H refined the protocol, ran
412 the experiments, collected the data, developed the software, and analyzed the data. J.K.M was the
413 primary author of the paper with significant contributions from S.K.H and A.v.d.B. All authors were
414 involved in the revision of the draft manuscript and have agreed to the final content.

415 COMPETING INTERESTS

416 The authors have no financial, personal, or professional competing interests that could be construed
417 to unduly influence the content of this article.

418 GRANT INFORMATION

419 The work was partially funded by the State of Ohio Third Frontier Commission through the Wright
420 Center for Sensor Systems Engineering (WCSSE).

421 ACKNOWLEDGMENTS

422 We thank Roman Boychuk and Obinna Nwanna for assistance with the experiments.

423 REFERENCES

- 424 Chester, V. L., Tingley, M., and Biden, E. N. (2007). Comparison of two normative paediatric gait
425 databases. *Dynamic Medicine*, 6:8.
- 426 Kirtley, C. (2014). CGA Normative Gait Database. <http://www.clinicalgaitanalysis.com/data/>.
- 427 Moore, J. K., Hnat, S., and van den Bogert, A. (2014a). Dynamic gait data collected during walking
428 under the influence of random perturbations on an actuated treadmill.
- 429 Moore, J. K., Nwanna, O., Hnat, S., and van den Bogert, A. (2014b). Gaitanalysis toolkit: Version
430 0.1.2.
- 431 Sup, F., Bohara, A., and Goldfarb, M. (2008). Design and control of a powered transfemoral
432 prosthesis. *The International Journal of Robotics Research*, 27(2):263–273.
- 433 Tirosh, O., Baker, R., and McGinley, J. (2010). GaitaBase: Web-based repository system for gait
434 analysis. *Computers in Biology and Medicine*, 40(2):201–207.

- 435 van den Bogert, A. J. (2003). Exotendons for assistance of human locomotion. *BioMedical*
436 *Engineering OnLine*, 2(1):17.
- 437 van den Bogert, A. J., Geijtenbeek, T., Even-Zohar, O., Steenbrink, F., and Hardin, E. C. (2013). A
438 real-time system for biomechanical analysis of human movement and muscle function. *Medical*
439 *& Biological Engineering & Computing*, pages 1–9.
- 440 Vaughan, C., Davis, B., and O'Connor, J. (1992). *Dynamics of Human Gait*. Human Kinetics
441 Publishers, 1st edition.
- 442 Wang, Y. and Srinivasan, M. (2014). Stepping in the direction of the fall: the next foot place-
443 ment can be predicted from current upper body state in steady-state walking. *Biology Letters*,
444 10(9):20140405.
- 445 White, E. P., Baldrige, E., Brym, Z. T., Locey, K. J., McGlenn, D. J., and Supp, S. R. (2013). Nine
446 simple ways to make it easier to (re)use your data. *PeerJ PrePrints*, 1:e7v2.
- 447 Willson, J. D. and Kernozek, T. (2014). Gait data collected at univ of wisconsin-LaCrosse.
- 448 Winter, A., D. (1990). *Biomechanics and Motor Control of Human Movement*. 2nd edition.

```

study:
  id: 1
  name: Gait Control Identification
  description: Perturb the subject during walking and running.
subject:
  id: 8
  age: 20
  mass: 70.0
  mass-units: kilograms
  height: 1.572
  height-units: meters
  knee-width-left: 107.43
  knee-width-right: 107.41
  knee-width-units: millimeters
  ankle-width-left: 70.52
  ankle-width-right: 67.66
  ankle-width-units: millimeters
  gender: male
trial:
  id: 58
  subject-id: 8
  datetime: 2014-03-28
  notes: >
    The subject did a somersault during this trial instead of following
    instructions to walk. Will have to use for another study.
  nominal-speed: 0.8
  nominal-speed-units: meters per second
  stationary-platform: True
  pitch: False
  sway: False
  hardware-settings:
    high-performance: True
  dflow-version: 3.16.1
  cortex-version: 3.1.1.1290
  marker-map:
    M1: LHEAD
    M2: THEAD
    M3: RHEAD
    M4: FHEAD
    M5: C7
  analog-channel-map:
    Channel1.Anlg: F1Y1
    Channel2.Anlg: F1Y2
    Channel3.Anlg: F1Y3
    Channel4.Anlg: F1X1
  events:
    A: Force Plate Zeroing
    B: Calibration Pose
    C: First Normal Walking
    D: Longitudinal Perturbation
    E: Second Normal Walking
    F: Unloaded End
  files:
    compensation: ../T057/mocap-057.txt
    mocap: mocap-058.txt
    record: record-058.txt
    meta: meta-058.yml

```

Listing 1. A fictitious example of a YAML formatted meta data file. All of the possible keys in the data set are shown.

```
>>> from gaitanalysis.motek import DFlowData
>>> data = DFlowData('mocap-020.txt', 'record-020.txt',
...                  'meta-020.yml')
>>> mass = data.meta['subject']['mass']
>>> data.clean_data()
>>> event_df = dflow_data.extract_processed_data(
...     event='Longitudinal Perturbation')
```

Listing 2. Python interpreter session showing how one could load a trial into memory, extract the subject's mass from the meta data, run the data cleaning process, and finally extract a Pandas DataFrame containing all of the time histories for a specific event in the trial.

```
>>> from gaitanalysis.gait import GaitData
>>> gdata = GaitData(event_df)
>>> gdata.inverse_dynamics_2d(left_markers, right_markers,
...                           left_loads, right_loads, mass, 6.0)
>>> gdata.grf_landmarks('Right Fy', 'Left Fy', threshold=20.0)
>>> gdata.split_at('right')
>>> gdata.plot_gait_cycles('Left Hip Joint Torque', mean=True)
>>> gdata.save('gait-data.h5')
```

Listing 3. Python interpreter session showing how one could use the GaitData class to load in the result of DFlowData and compute the inverse dynamics (joint angles and torques), identify the gait event (e.g. heel strikes), split the data with respect to the gait events in a Pandas Panel, plot the mean and standard deviation of one time history with respect to the gait cycles, and save the data to disk.