

We would like to thank the reviewer again for taking the time to provide further feedback, which has been constructive and valuable in improving this manuscript. We have addressed the additional points outlined by the reviewer: in particular we have amended the text in Sections 3.1 and 3.2 of the Results to give a greater commentary on the role of model structure and ecosystem behavior in the explanation of our results.

On a strictly factual basis, the authors have responded to my initial comments, and made changes to the paper that address them. It still seems to me that the discussion is fairly clinical, especially for Figures 2-4. I think that there is more information here than what the authors are telling us, that what the reader is provided is more basic reportage than insight into savanna behavior and models' representation of it. I'm sympathetic to the fact that there is a massive amount of information on these graphs, and not every detail can be described. But it seems to me that there are causal explanations, having to do with model structure and ecosystem behavior, beyond a recitation of which models did what.

But I appreciate that the authors are showing us model results in Section 3, and then talking about implications and interpretation in the Discussion and conclusions (Sections 4 and 5).

I do still have a few comments and questions about figures 3 and 4:

Comment [Figure 3]:

There is a paragraph about the empirical benchmarks (lines 356-363), but not really in the context of what they show in this figure. There are places where emp2/emp3 perform quite well (dry season LE Howard Spgs), places where none of them seem to have a handle (dry season GPP, Dry River). The authors briefly mention the lack of internal modes and storage capability in emp1-3, but we're not given much more than that. For example, when emp1 has a peak at a similar value to what is observed, the density is generally much too high. When the peak density

is reasonable, there is generally an offset in value. What does this mean? Anything? Does emp1 have value at all other than a low bar for models to clear?

Response 1: The purpose of the benchmarking in this study was to give an indication as to how well we should expect a model to perform given its level of complexity. It is important to remember that the benchmarks are simple regression models, trained on the meteorological forcing (e.g. $X_{1,t}$, $X_{2,t}$); they are not TBMs. Consequently the benchmarks (emp1...3) used here are not competing models per se, and their purpose is not to elucidate any particular behaviours of the ecosystem. Rather they provide an extra point of reference by which model performance can be quantified. For example, a TBM that might be driven purely by solar irradiance (R_s) would be expected to perform no better or worse than the emp1 benchmark used in this study. It would be unrealistic to expect such a model to have the same predictive capability to simulate LE or GPP to within a low degree of error; there just isn't enough information provided by R_s alone for that comparison to be fair. Assessing model performance in this way, therefore allows us to quantify the predictive capability of a TBM by both the minimisation of model error (i.e. model-data comparison) and the degree of model complexity needed to achieve it.

We have edited the text and added additional information to Section 3.1 in order to: i) give greater reference to the role of model structure and ecosystem behaviour in the explanation of results (meeting the request in the reviewer's opening comment), and ii) give greater clarity to the role of the empirical benchmarks. The additional text is quoted below:

Figure 3 shows the probability density functions (PDFs) for the wet (Nov – Apr) and dry season (May – Oct) fluxes at each site. Tower and model PDFs were determined by binning each flux into the respective seasons and using kernel density estimation (Bashtannyk and Hyndman, 2001) to determine smoothed distributions. The shape and mean position of the distributions indicate the ability of the models to capture the extremes (day-to-day variability) and the seasonality of the fluxes respectively, highlighting possible predictive biases (i.e. the over- or underestimation of the tower fluxes). Across the NATT, the PDFs for the tower fluxes tended to shift to low values and became narrower as annual rainfall declined, and this was most prominent in the dry

season. A change in the spread and mean position of the flux tower PDFs demonstrate the strong seasonality of water and carbon exchange at all sites. By contrast, the PDFs of the model simulations did not replicate this trend, having high densities and being mostly stationary across sites. Regarding savanna water-use, the distributions of the BIOS2 and SPA models were similar to those of the flux towers. The BESS model also showed a similar distribution of LE, despite the fact that it did not simulate soil water extraction. The LPJGUESS model, which had the shallowest simulated tree rooting depth, displayed PDFs of high density that were biased towards low LE ($20 - 40 \text{ W m}^{-2}$) across all sites and seasons. The MAESPA model showed a similar behaviour, despite this model having a much deeper simulated rooting depth and a root-water extraction scheme that is equivalent to the SPA model. The distributions for the CABLE and BIOS2 models were largely disparate despite these models being functionally equivalent. Notably, CABLE wet season LE was more broadly distributed ($5 - 200 \text{ W m}^{-2}$) than the flux towers and other models at all sites, while dry season LE was narrower. In relation to savanna carbon uptake, all models showed wet and dry season PDFs of high density that became more closely aligned with the flux tower distributions as the sites became drier. The behaviour of the modelled GPP distributions were otherwise similar to those of the modelled LE distributions. The differences among TBM and flux tower PDFs indicated possible issues in simulated processes that are active during the wet season.

The benchmarks set low to high levels of expected TBM performance across the NATT. Additionally, they also demonstrated the level of model complexity that is required to simulate water and carbon exchange at these sites. The simplest of the benchmarks, represented as a linear regression of the fluxes against R_s (emp1), was capable of predicting the magnitude and daily time-course of the tower fluxes (data not shown), but there was not enough information in R_s to capture the seasonality or the distribution of the fluxes expressed by the tower data. The intermediate benchmark that included additional meteorological information on T_a and VPD (emp2) gave an improved capability in capturing the flux distributions, but could not replicate the full seasonality of the fluxes across the NATT. It was only by including additional phenological information (LAI) together with site meteorology (R_s , T_a and VPD) that the seasonality and distribution of the fluxes could be captured, as demonstrated by the most complex benchmark (emp3). This indicated that in order for the TBMs to achieve the best possible performance at simulating water and carbon exchange along the NATT, the correct implementation and utilisation of phenological information by the models was required. All TBMs used in this study utilised this breadth of information, but only some

of the models were capable of meeting the expected level of performance set by the emp3 benchmark, and only then for specific sites and seasons.

Comment [Figure 4]:

Do the bias vs. flux plots add information beyond what is in the time-course plots? The other reviewer touched on this when asking about the regression lines-do they give us pertinent information or not? There is no discussion about what the bias vs. flux plots mean, beyond what the reader can see for him/herself. CABLE has positive slope at all sites for LE, and crosses zero. To me, this says that during the dry season, at all sites, CABLE underestimates LE. In the wet season, LE is overestimated. But we can see that from the time-course plot. In this case the bias/flux plot seems superfluous: Are there other cases where they are important to the analysis, or can they be discarded?

Response 2: The intention of the bias subplots was to give a clearer representation of the degree of over- and underestimation of predicted water and carbon exchange by the TBMs along the NATT, with the focus being on underestimation of fluxes in the dry season. However, the scatter of points in the residuals vs. model prediction subplots is not necessarily linear, and the use of a regression line is likely to mislead the reader to believe otherwise (as the reviewer points out here). We agree that using these subplots purely as a way of qualitatively assessing model bias does not add much value to the analysis and is likely to cause confusion in the interpretation of the results. They have therefore been discarded.

Comment [Figure 4 continued]:

The positive residuals for BIOS2 at Howard Springs and Adelaide River in Apr-Jun are interesting. That says to me that at these wetter sites, the 'real world' dries out quicker than the model does, and water is available for ET in the model. This doesn't happen in this model at the drier sites. Also, BIOS; Howard/Adelaide are negatively biased in the wet season, more than other sites. With a 10m deep soil for

water storage, and 5m root depth for access, I might think that BIOS would do a better job in these wetter sites, but instead the comparison is better in dry sites. Some commentary about some of these types of behaviors in the context of model architecture would be welcome, but I don't demand it.

Response: We have amended the entirety of Section 3.2 to reflect a greater commentary on how model structure may be affecting the behaviour of residual error. The text is quoted below:

An analysis of the model residuals was conducted to show how model structure affects the prediction of savanna fluxes across the rainfall gradient. To do this we examined the standardised model residuals from each TBM, determined by expressing the residual error in terms of its standard deviation. Figure 4 shows the residual time-series for model predicted LE and GPP at each savanna site and provides an effective way of examining how a model responds to progressive changes in the environment, through the expression of model bias and error (Medlyn et al., 2005).

The model residuals demonstrated that there was significant bias and heteroscedasticity in predicted LE and GPP in almost all cases. The residual time-series showed that model error was largest in the wet season, but declined with the transition into the dry season. Additionally, the models underestimated LE and GPP more significantly during the wet season. A possible explanation for this behaviour is that during the wet season, multiple land-surface components: the soil surface, the understorey grasses, and the tree canopy (i.e. 3 sources for potential error) contribute to the bulk fluxes, while during the dry season only the tree canopy contributes (i.e. 1 source for potential error). It is likely that the reduction in residual error between wet and dry seasons was a result of the declining influence of the grasses and the soil surface to ecosystem land-surface exchange during the latter period (via senescence and low surface soil moisture respectively). The bias towards the underestimation of wet season fluxes was more pronounced at the mesic sites (Howard Springs, Adelaide River), despite some models simulating relatively deep root profiles (e.g. BIOS2, MAESPA). Differences in how the TBMs simulated root-water extraction also had no effect on reducing this bias (e.g. MAESPA, SPA). Given that soil-water was not a limiting factor at the mesic sites during this period, deep root profiles offered limited advantage towards model performance. Nonetheless, the simulated tree root-zone appeared to be an important factor for all sites during the dry season, with shallow root depths (LPJGUESS: 2 m) and/or

inadequate root-water uptake schemes (CABLE: concentrated in the upper soil profile) the likely cause for underestimation during this period. However, as the sites became drier (e.g. Sturt Plains) a shallow root-profile was suitable to give flux estimates of a reasonably low error. Despite model error reducing with the increase in ecosystem water limitation that occurs in both space (down the NATT) and time (wet to dry season), there are still patterns of model bias that may be unrelated to simulated soil-water dynamics. This is particularly obvious during the wet-to-dry transition periods (e.g. BIOS2, SPA) when the C₄ grass understorey senesces, indicating possible problems with the how the models translate information on phenology.

Comment [Figure 4 continued]:

The bias vs. flux for LPJGUESS/GPP show a negative slope for all sites but the driest two, which have a positive slope. Does this mean anything, or are these slopes indicating trends that really aren't there (it appears that the bias at Sturt Plains is always positive). If the latter, is it possible that these regression lines are suggesting facts not in evidence? If there is actually 'a larger tendency towards underestimation in the wet than in the dry season' then we should see a negative bias at larger GPP/LE at all sites. I see this for CABLE;GPP, LPJGUESS;LE/GPP, MAESPA;LE. But I don't see it for other metrics (CABLE; LE, BESS;LE/GPP, MAESPA; GPP). Am I misinterpreting something?

Response 4: We refer the reviewer to response 2

Typos and errata

Line 100: change 'savannas' to 'savanna'

Response: Done.

Line 126: change 'their ability at simulating' to 'their ability to simulate'

Response: Done.

Line 143: change 'savannas' to 'savanna'

Response: Done.

Line 169: change 'Table 1, and' to 'Table 1, which'

Response: Done.

Line 174: Period is on the subscript

Response: Done.

Line 212: change 'simulate water; carbon' to 'simulate water and carbon'

Response: Done.

Line 220: change 'BIOS is fine-spatial-resolution' to 'BIOS is a fine spatial resolution'

Response: Done.

Line 223: change 'similarities reduces' to 'similarities reduce'

Response: Done.

Line 223: remove comma after 'truly'

Response: Done.

Line 260: delete 'then'

Response: Done.

Line 261: PALS is given as acronym before it is described in lines 266-267

Response: Done.

Line 308: should be 'Table A1', not 'Table B1'

Response: Done.

Line 319: delete 'in the determination of these metrics'

Response: Done.

Line 329: add 'each' after 'flux and'

Response: Done.

Line 339: delete 'savanna'. All sites are savanna sites

Response: Done.

Line 379: delete 'vegetation'. Transpiration is from vegetation by definition.

Response: Done.

Line 382: delete 'and' (first word of line)

Response: Done.

Line 457: insert 'it' between 'although' and 'performed'

Response: Done.

Line 495: should be 'shrubs are established' or 'shrubs have been established'

Response: Done.

Lines 530-533: once you say that up to 75% of total LE is from understory herbaceous transpiration, mentioning that C4 understory will be a large contributor is redundant.

Response: We have decided to keep the structure of the sentence, with a change of phrase from '*C4 understory*' to '*C4 grasses*'. We believe the statement makes explicit that soil evaporation contributes very little to the water flux occurring in the understory.

Line 560: change 'of TBMs' to 'of the TBMs'

Response: Done.

Line 561: delete 'that we discuss below'

Response: Done.

Line 592: change semi-colon at the end of the line to a comma

Response: Done.

Lines 622, 623: change semi-colons to commas

Response: Done.

Line 654: delete 'on'

Response: Done.

Line 659. Remove dash before 'The'

Response: Done.

1 **A model inter-comparison study to examine limiting factors in**
2 **modelling Australian tropical savannas**

3

4 | **Authors:** Rhys Whitley¹, Jason Beringer², Lindsay B. Hutley³, Gab Abramowitz⁴, Martin
5 G. De Kauwe¹, Remko Duursma⁵, Bradley Evans⁶, Vanessa Haverd⁷, Longhui Li⁸,
6 Youngryel Ryu⁹, Benjamin Smith¹⁰, Ying-Ping Wang¹¹, Mathew Williams¹², Qiang Yu⁷

7

8 **Institutions:**

9 ¹Department of Biological Sciences, Macquarie University, North Ryde, NSW 2109,
10 Australia

11 ²School of Earth and Environment, University of Western Australia, Crawley, WA 6009,
12 Australia

13 ³School of Environment, Charles Darwin University, Casuarina, NT 0810, Australia

14 ⁴Climate Change Research Centre, University of New South Wales, Kensington, NSW
15 2033, Australia

16 ⁵Hawkesbury Institute for the Environment, University of Western Sydney, Penrith, New
17 South Wales 2751, Australia

18 ⁶Faculty of Agriculture and Environment, University of Sydney, Eveleigh, NSW 2015,
19 Australia

20 ⁷CSIRO Ocean and Atmosphere, Canberra 2601, Australia

21 ⁸School of Life Sciences, University of Technology Sydney, Ultimo, NSW 2007, Australia

22 ⁹Department of Landscape Architecture and Rural Systems Engineering, Seoul National
23 University, Seoul, South Korea

24 ¹⁰Department of Physical Geography and Ecosystem Science, Lund University, Lund,
25 Sweden

26 ¹¹CSIRO Ocean and Atmosphere, Aspendale, Victoria 3195, Australia

27 ¹²School of GeoSciences, University of Edinburgh, Edinburgh, United Kingdom

28

29 **Corresponding author:**

30 Name: Rhys Whitley

31 Email: rhys.whitley@mq.edu.au

32 Address:

33 Department of Biological Sciences

34 Macquarie University

35 North Ryde, NSW

36 2109, Australia

37

38

39 **Abstract:**

40 Savanna ecosystems are one of the most dominant and complex terrestrial biomes that
41 derives from a distinct vegetative surface comprised of co-dominant tree and grass
42 populations. While these two vegetation types co-exist functionally, demographically
43 they are not static, but are dynamically changing in response to environmental forces
44 such as annual fire events and rainfall variability. Modelling savanna environments with
45 the current generation of terrestrial biosphere models (TBMs) has presented many
46 problems, particularly describing fire frequency and intensity, phenology, leaf
47 biochemistry of C₃ and C₄ photosynthesis vegetation, and root water uptake. In order to
48 better understand why TBMs perform so poorly in savannas, we conducted a model
49 inter-comparison of 6 TBMs and assessed their performance at simulating latent energy
50 (LE) and gross primary productivity (GPP) for five savanna sites along a rainfall gradient
51 in northern Australia. Performance in predicting LE and GPP was measured using an
52 empirical benchmarking system, which ranks models by their ability to utilise
53 meteorological driving information to predict the fluxes. On average, the TBMs
54 performed as well as a multi-linear regression of the fluxes against solar radiation,
55 temperature and vapour pressure deficit, but were outperformed by a more complicated
56 nonlinear response model that also included the leaf area index (LAI). This identified
57 that the TBMs are not fully utilising their input information effectively in determining
58 savanna LE and GPP, and highlights that savanna dynamics cannot be calibrated into
59 models and that there are problems in underlying model processes. We identified key
60 weaknesses in a model's ability to simulate savanna fluxes and their seasonal variation,
61 related to the representation of vegetation by the models and root water uptake. We
62 underline these weaknesses in terms of three critical areas for development. First,
63 prescribed tree-rooting depths must be deep enough, enabling the extraction of deep
64 soil water stores to maintain photosynthesis and transpiration during the dry season.
65 Second, models must treat grasses as a co-dominant interface for water and carbon
66 exchange, rather than a secondary one to trees. Third, models need a dynamic
67 representation of LAI that encompasses the dynamic phenology of savanna vegetation
68 and its response to rainfall interannual variability. We believe this study is the first to
69 assess how well TBMs simulate savanna ecosystems, and that these results will be used
70 to improve the representation of savannas ecosystems in future global climate model
71 studies.

72

73 **Introduction**

74 Savanna ecosystems are a diverse and important biome that play a significant role in
75 global land-surface processes (van der Werf et al., 2008). Globally, they occupy regions
76 around the wet-dry tropical to sub-tropical equatorial zone, covering approximately 15
77 to 20% of the terrestrial surface and contribute ~30% to global net primary production
78 (Grace et al., 2006; Lehmann et al., 2014). Savannas are water-limited ecosystems where
79 rainfall is often seasonal or monsoonal, and have a spatial extent that can cover an area
80 with annual rainfall in the range of 500 to 2000 mm (Bond, 2008; Kanniah et al., 2010;
81 Sankaran et al., 2005). The variability in the amount and timing of annual rainfall,
82 coupled with local topo-edaphic properties, and the frequency and intensity of seasonal
83 fires strongly influences the structure and function of savanna vegetation (Beringer et
84 al., 2007; Kanniah et al., 2010; Ma et al., 2013; Sankaran et al., 2005). Savannas are
85 characterised by a multi-layer stratum of vegetation, where an open and discontinuous
86 canopy overstorey is seasonally dominated by understorey grasses (Scholes and Archer,
87 1997). These tree and grass layers are distinctly and functionally different, fixing carbon
88 using different photosynthetic pathways, C₃ and C₄ photosynthesis respectively (Bond,
89 2008; Scholes and Archer, 1997; Williams et al., 1996b). The canopy overstorey can be
90 either evergreen or deciduous (depending on the evolutionary history), while the grass
91 understorey is annual: active only in the wet season and senescing at the end of this
92 period (Williams et al., 1996b). Consequently, water, carbon and nutrient cycling in
93 savannas is largely determined from the balance and co-existence of these two life forms
94 (Lehmann et al., 2009; Sankaran et al., 2005).

95 Given the complex nature of savannas, modelling the land surface exchange and
96 vegetation dynamics for this biome is challenging for terrestrial biosphere models
97 (TBMs). Here we define TBMs to broadly encompass stand, land-surface, and dynamic
98 global vegetation models (Pitman, 2003). Most land surface schemes that feed into
99 larger earth system models use simplistic representations of vegetation, and these will
100 have difficulty describing the complex structure of savanna ecosystems. Such issues may
101 be: simplistic assumptions in relation to rooting depth and inadequate responses to
102 drought (De Kauwe et al., 2015; Li et al., 2012); ignoring the multilayered nature of
103 savannas and the differing structural (including radiation), functional (including
104 different plant functional types) and phenological differences (Whitley et al., 2011); and
105 in some cases neglecting the C₄ photosynthetic pathway entirely (Parton et al., 1983;
106 Schymanski et al., 2007) It is therefore critical that TBMs meet the challenges that

Rhys Whitley 10/5/16 12:07 PM

Deleted: s

108 savanna dynamics present if water and carbon exchange are to be correctly simulated in
109 response to global change.

110 Despite these issues, there have been significant advances in modelling savanna
111 dynamics in recent years, and this has been focused on integrating important features
112 specific to savanna ecosystems, namely frequent fire and tree-grass competitive
113 interactions, processes that shape savanna structure and function (Haverd et al., 2016;
114 Higgins and Scheiter, 2012; Scheiter and Higgins, 2007; Scheiter et al., 2014; Simioni et
115 al., 2003). Nevertheless, little work has been undertaken to critically evaluate the
116 performance and processes of TBMs when used to capture water and carbon cycling in
117 savannas, most notably in west Africa (Simioni et al., 2000) and Australia (Schymanski
118 et al., 2007, 2008, 2009; Whitley et al., 2011). Many global ecosystem models moreover
119 use broad plant functional types (PFTs) with single parameter values to describe whole
120 biomes (Pitman, 2003), making them unable to represent changing vegetation structure
121 (tree:grass ratio) in the continuum of grassland to woodland savanna. Approaches have
122 been developed that can account for savanna dynamics, such as using mixed tiles,
123 whereby trees and grasses are simulated as separate surfaces that are then aggregated
124 together (Kowalczyk et al., 2006). However, this approach fails to capture the
125 competition between trees and grasses for light, water and nutrient resources.

126 In this study, we take 6 TBMs of distinctly different conceptual frameworks, and assess
127 their ability to simulate savanna water and carbon exchange along the North Australian
128 Tropical Transect (NATT) that is defined by a strong rainfall gradient. Australian
129 tropical savannas can be considered largely intact compared to South American and
130 African savannas, and provide a 'living laboratory' to understand the links between
131 vegetation structure and function and how it responds to environmental change (Hutley
132 et al., 2011). We challenge the models by evaluating them along the rainfall gradient,
133 which extends over a broad biogeographical extent and strong interannual variability in
134 climate (Koch et al., 1995). The aim of this study is to highlight critical processes that
135 may be missing in current TBMs and are required to adequately simulate savanna
136 ecosystems. Specifically, we examine whether a TBM's structural framework, such as the
137 representation of the understorey grasses (C₄ photosynthesis), tree rooting depth, and
138 description of phenology (prescribed vs. dynamic) can adequately replicate observed
139 carbon and water fluxes. To achieve this we measure the performance of each TBM by
140 comparing its predictions to a set of empirical benchmarks that describe *a priori*
141 expected levels of model performance. We identify regions of low performance among
142 sites and seasons, to diagnose under what climate conditions reduced model

Rhys Whitley 10/5/16 12:08 PM

Deleted: at

Rhys Whitley 10/5/16 12:08 PM

Deleted: ing

145 performance occurs. We then infer what processes (present or missing) may be the
146 cause for reduced performance when applied to savanna ecosystems. Our intention is
147 that these results can be used to flag high priorities for future development by the
148 terrestrial biosphere modelling community.

149

150 2. Methodology

151 2.1 Observational data

152 The North Australian Tropical Transect (NATT) is a sub-continental rainfall gradient in
153 the wet-dry tropical climate zone of Northern Australia, which encompasses a distance
154 of approximately 1000 km over a latitudinal range of -12 to -23 °S and a decline in mean
155 annual precipitation (MAP) from 1700 mm to 300 mm (Hutley et al., 2011). It is one of
156 three savanna transects established in the mid 1990's, forming part of the International
157 Geosphere Biosphere Program (IGBP) along with the SAVannas in the Long Term (SALT)
158 transect in West Africa and the Kalahari Transect (KT) in Southern Africa (Koch et al.,
159 1995). Soils range from sand dominated red Kandosols to black, cracking clay soils that
160 are more extensive in the southern end of the NATT that are limiting to woody plant
161 growth (Hutley et al., 2011; Williams et al., 1996b). Kandosols are ancient and
162 weathered, such that they have been leached of nutrients by the large monsoonal
163 rainfall (McKenzie et al., 2004). Close to the northern coastline, vegetation is comprised
164 primarily of evergreen *Eucalyptus* and *Corymbia* tree species that overlie an understorey
165 of C₄ *Sorghum* and *Heteropogon* spp. grasses. Inland, tree biomass, leaf area index (LAI)
166 and cover tends to decline and by -18 °S savanna vegetation transitions to less dense
167 *Acacia* woodlands, shrublands and grasslands that are dominated by *Astrelba* grass
168 species (Hutley et al., 2011). Fires occur regularly in these environments, increasing in
169 frequency with higher rainfall (MAP > 1000 mm), and are fuelled by the accumulation of
170 understorey C₄ grasses that cure in the dry season (Beringer et al., 2014; Russell-Smith
171 and Edwards, 2006).

172 The five flux tower sites along the NATT used in this study are outlined in Table 1, which
173 describes stand soil and vegetation characteristics, as well as a summary of local
174 meteorology (Hutley et al., 2011). These sites represent a sampling of savanna
175 environments covering a wide range of MAP and a much smaller range of mean annual
176 temperature (MAT) (Fig. 1). At each site, an eddy covariance system was used to
177 measure the ecosystem-atmosphere exchange of radiation, heat, water and CO₂ Quality

Rhys Whitley 10/5/16 12:09 PM

Deleted: s

Rhys Whitley 10/5/16 12:09 PM

Deleted: and

Rhys Whitley 10/5/16 12:10 PM

Formatted: Not Superscript/ Subscript

180 assurance and control (QA/QC) and corrections on the fluxes were carried out on the 30
181 minute dataset using the OzFlux QC/QA protocol (v2.8.5), developed by the OzFlux
182 community under creative commons licensing (www.ozflux.org.au) (see Eamus et al.,
183 2013). Missing or rejected data were gap-filled using the DINGO (Dynamic INtegrated
184 Gap filling and partitioning for Ozflux) system (see Moore et al., 2016). Gross primary
185 productivity (GPP) was not observed but determined from the difference between
186 measured net ecosystem exchange (NEE) and modelled ecosystem respiration (Re).
187 Values of Re were determined by assuming nocturnal NEE equals Re under the
188 conditions for sufficient turbulent transport. Values that meet these requirements are
189 then used to make daytime predictions of Re, using an artificial neural network (ANN),
190 with soil moisture and temperature, air temperature, and the normalised difference
191 vegetation index (NDVI) used as predictors. Additionally, the effect of fire on the water
192 and carbon fluxes are quantified and incorporated into the datasets accounting for the
193 nonlinear response in productivity (becoming a carbon source) during the post-fire
194 recovery period (Beringer et al., 2007). Because the TBMs used here do not attempt to
195 simulate stochastic fire events (and other disturbance regimes), these post-fire recovery
196 periods were removed when determining the benchmarks and model performance as
197 described below.

198 Finally, we use the definitions for water and carbon exchange as outlined by Chapin et
199 al. (2006), whereby the sub-daily rate of GPP is expressed in $\mu\text{mol m}^{-2} \text{s}^{-1}$ and uses a
200 negative sign (-) to denote the removal of CO_2 from the atmosphere. Similarly, LE is
201 expressed in terms of energy as W m^{-2} and uses a positive sign to denote the addition of
202 H_2O to the atmosphere.

203

204 **2.2 Terrestrial biosphere models**

205 The 6 TBMs used in this study cover a wide spectrum of characteristics of operation,
206 scale and function, and include differences in operational time-step (30min vs. daily),
207 scope of simulated processes (soil hydrology, static or dynamic vegetation, multi-layer
208 or big leaf description of the canopy) and intended operational use (coupled to earth
209 system models, offline prediction, driven by remote sensing products). These
210 characteristics along with what we define as a model 'functional class' are given in Table
211 2 and are defined as follows. Stand models (SMs) give detailed multi-layer descriptions
212 of canopy and soil processes for a particular point, operating at a sub-daily time-step

213 (Soil-Plant-Atmosphere model: SPA, and MAESPA). Land-Surface models (LSMs) operate
214 at the same temporal resolution as SMs, but adopt a simpler representation of canopy
215 processes, allowing them to be applied spatially (Community Atmosphere Biosphere
216 Land Exchange model; CABLE, and BIOS2; a modified version of CABLE). Dynamic Global
217 Vegetation Models (DGVMs) simulate water and carbon much like the other models, but
218 simulate dynamic rather than static vegetation that changes in response to climate and
219 disturbance (Lund-Potsdam-Jena General Ecosystem Simulator; LPJGUESS). Lastly,
220 Remote Sensing models (RSMs) are driven by remotely sensed atmospheric products,
221 and infer water-stress of vegetation through changes in fractional cover rather than
222 detailed soil hydrological processes (Breathing Earth System Simulator; BESS). Some of
223 the TBMs share similar structural frameworks in parts: for example, both SPA and
224 MAESPA use similar below-ground soil hydrology and root-water uptake schemes, while
225 BIOS2 is fine spatial resolution (0.05 degree), offline modelling environment for
226 Australia, in which predictions of CABLE (with alternate parameterisations of drought
227 response and soil hydrology) are constrained by multiple observation types (see Haverd
228 *et al.* 2013). Although these similarities reduce the number of truly functionally,
229 independent models used in the experiment, the presence of such overlap can be useful
230 in identifying if particular frameworks are the cause for model success or failure.

231

232 2.3 Experimental protocol

233 All TBMs were parameterised for each of the five savanna sites using standardised
234 information on vegetation and soil profile characteristics (Table 1). For TBMs that
235 required them, parameter values pertaining to leaf biochemistry, such as maximum
236 Rubisco activity (V_{cmax}) and leaf nitrogen content per leaf area (N_{area}), were assigned
237 from Cernusak *et al.* (2011), who undertook a physiological measurement campaign
238 during the SPECIAL program (Beringer *et al.* 2011). Parameters relating to soil sand and
239 clay content were taken from the Australian Soil Classification (Isbell, 2002), while root
240 profile information was sourced from Chen *et al.* (2003) and Eamus *et al.* (2002). Each
241 TBM was setup to describe a C₃ evergreen overstorey with an underlying C₄ grass
242 understorey, and conforms well with the characteristics of savannas in Northern
243 Australia (Bowman and Prior, 2005). All TBMs (excluding LPJGUESS) prescribed LAI as
244 an input, to characterise the phenology of vegetation at each site. In these cases LAI was
245 determined from MODIS derived approximations that were well matched to ground-
246 based estimations of LAI at the SPECIAL sites (Sea *et al.*, 2011). The fraction of C₃ to C₄

Rhys Whitley 10/5/16 12:13 PM
Deleted: ,

Rhys Whitley 10/5/16 12:14 PM
Deleted: -

Rhys Whitley 10/5/16 12:14 PM
Deleted: -

Rhys Whitley 10/5/16 12:14 PM
Deleted: s

Rhys Whitley 10/5/16 12:14 PM
Deleted: ,

252 vegetation was handled differently by each model and was determined for each as
253 follows. For MAESPA and SPA, the models allowed for time-varying tree and grass
254 fractions to be assigned as direct inputs, and these time-varying fractions were
255 determined using the method of Donohue et al. (Donohue et al., 2009). BIOS2 similarly
256 used the same method to extract time-varying fractions, while CABLE used a static
257 fraction that did not change. The BESS model derived the C₃:C₄ fraction from the C₃ and
258 C₄ distribution map of Still et al. (2003), while for LPJGUESS this fraction is a prognostic
259 determination resulting from the competition between trees and grasses (see Smith et
260 al., 2001). Model simulations were driven using observations of solar radiation, air
261 temperature, relative humidity (or vapour pressure deficit; VPD), rainfall, atmospheric
262 CO₂ concentration and LAI (if prescribed), and included a spin-up period of 5 years to
263 allow internal states, such as the soil water balance and soil temperature to reach
264 equilibrium. The exception to the above was the BIOS2 model, which was run using
265 gridded meteorological inputs and had its model parameters optimised through a
266 model-data fusion process (see Haverd et al., 2013).

267 Simulations for each savanna site covered a period of 2 to 10 years depending on the
268 availability of data from each flux site (Table 1) and results were standardised to the
269 ALMA (Assistance for Land-surface Modelling Activities) convention. Model predictions
270 of LE and GPP were evaluated against local observations at each site from the eddy
271 covariance datasets and benchmarked following the methodology proposed by the
272 [Protocol for the Analysis of Land-surface models \(PALS\) and the PALS](#) Land Surface
273 Model Benchmarking Evaluation Project (PLUMBER) (Abramowitz, 2012; Best et al.,
274 2015) as described below.

275

276 2.4 Empirical benchmarking

277 The paradigm for model assessment [outlined by PALS](#), (Abramowitz, 2012) suggests that
278 model assessment is more meaningful when *a priori* expectations of performance in any
279 given metric can be defined. Such benchmarks can be created using simple empirical
280 models, built on statistical relationships between the fluxes and drivers, and establish
281 the degree to which models utilise the information available in their driving data about
282 the fluxes they aim to predict. Additionally, these empirical models are simple in the
283 sense that they are purely instantaneous response to time-varying meteorological
284 forcing and contain no internal states or expression of ecophysiological processes. This

Rhys Whitley 10/5/16 12:15 PM

Deleted: then

Rhys Whitley 10/5/16 12:15 PM

Deleted: PALS

Rhys Whitley 10/5/16 12:17 PM

Deleted: set out in

Rhys Whitley 10/5/16 12:16 PM

Deleted: the Protocol for the Analysis of Land-surface models (

Rhys Whitley 10/5/16 12:16 PM

Deleted:)

291 is in comparison to TBMs that are complex, having some 20+ soil and vegetation
292 parameters, internal states, partitioning of light, as well as soil and vegetation, carbon
293 and nitrogen pools (Abramowitz et al., 2008).

294 We created a set of 3 empirical models of increasing complexity following the procedure
295 of Abramowitz (2012), which we compared with the TBMs. The first benchmark (emp1)
296 is simply a linear relationship between a turbulent flux (LE or GPP) and downward
297 short-wave radiation (R_s). The second benchmark (emp2) is slightly more complex, and
298 is a multi-linear regression between a flux and R_s , air temperature (T_a), and vapour
299 pressure deficit (VPD). Finally, the third benchmark (emp3) is the most complex and is a
300 nonlinear regression of the fluxes against R_s , T_a , VPD and LAI, determined from an ANN.
301 This benchmark is constructed using a self-organising linear output map that clusters
302 the four covariates into 10^2 distinct nodes and performs a multi-linear regression
303 between the fluxes and the 4 covariates at each node, resulting in a nonlinear (piece-
304 wise linear) response to the meteorological forcing data (Abramowitz et al., 2008; Hsu,
305 2002). In a departure from Abramowitz (2012), we include LAI as an additional
306 covariate, as the seasonal variance of savanna water and carbon exchange is strongly
307 coupled to the phenology of the grasses and to the deciduous and semi-deciduous
308 woody species (Moore et al., 2016). The seasonal behaviour of the empirical benchmark
309 drivers along the NATT can be referred to in the supplementary information. Empirical
310 benchmarks are created for each of the five flux sites using non gap-filled data, and are
311 parameterised *out-of-sample*, such that they use data from all sites except the one in
312 question. For example, the Howard Springs empirical benchmark models would use
313 information from Adelaide River, Daly Uncleared, Dry River and Sturt Plains to establish
314 their parameter values, but would exclude Howard Springs itself. Constructing the
315 benchmarks out-of-sample results in what is effectively a generalised response to an
316 independent dataset. Once the empirical models were calibrated for each site,
317 benchmarks were then created for both fluxes using the same meteorological forcing
318 used to run the TBMs.

319 Finally, we assess ecosystem model performance in terms of a ranking system, following
320 the PLUMBER methodology of Best et al. (2015). The performance of each individual
321 ecosystem model in predicting both LE and GPP at each site was determined using four
322 statistical metrics that describe the mean and variability of a model compared to the
323 observations. These metrics included the correlation coefficient (r), standard deviation
324 (sd), normalised mean error (NME), and mean bias error (MBE) (see Table A1).
325 Similarly, the same metrics were determined for each of the 3 benchmarks at each

Rhys Whitley 10/5/16 12:18 PM

Deleted: (Moore et al., *this issue*)

Rhys Whitley 10/5/16 12:20 PM

Deleted: B

328 savanna site. Each TBM was then ranked against the benchmarks (independently of the
329 other models) for each of the metrics listed above., where the ranking is between 1 and
330 4 (1 model + 3 benchmarks) and the best performing model for a given metric is ranked
331 as 1. An average ranking is then determined across all metrics for each TBM and all
332 benchmarks to give a final ranking of performance for each savanna site. The ranks
333 denote the number of metrics being met by the models and are not a measure of the
334 smallest absolute error. In determining the average ranks, the metrics were evaluated at
335 the daily time scale, as this was the lowest temporal resolution common amongst the 6
336 TBMs. Additionally, days where either driver or flux had been gap-filled were removed.
337 Herewith we use the term *performance* to relate to how well the TBMs compare to the
338 benchmarks as expressed by the ranks.

339

340 3. Results

341 3.1 Model predictions

342 Figure 2 shows the daily time-course of LE and GPP from the flux tower, models, and
343 benchmarks at each of the five savanna sites. Models, benchmarks and observations are
344 represented as a smoothed time-series (7-day running mean) and have been aggregated
345 into an ensemble year to express the typical seasonality of savanna water and carbon
346 exchange. Visually, the TBMs showed varying levels of performance across the rainfall
347 gradient. None of the models showed a clear consistency in simulating either flux, and
348 each responded differently to the meteorological drivers across sites. Additionally, some
349 of the models, such as CABLE and LPJGUESS, showed difficulty in simulating the
350 seasonality of the fluxes across the transect, particularly GPP. Differences among model
351 simulated LE and GPP were larger in the wet season than the dry season. However,
352 modelled LE and GPP appeared to co-vary quite strongly; overall both fluxes were
353 underestimated across sites by most models. Simulations by SPA and MAESPA were the
354 exception to this, broadly capturing tower GPP despite consistently underestimating LE
355 across sites.

356 Figure 3 shows the probability density functions (PDFs) for the wet (Nov – Apr) and dry
357 season (May – Oct) fluxes at each site. Tower and model PDFs were determined by
358 binning each flux into the respective seasons and using kernel density estimation
359 (Bashtannyk and Hyndman, 2001) to determine smoothed distributions. The shape and
360 mean position of the distributions indicate the ability of the models to capture the

Rhys Whitley 10/5/16 12:20 PM

Deleted: in the determination of these metrics

Rhys Whitley 10/5/16 12:22 PM

Deleted: savanna

Rhys Whitley 11/5/16 2:52 PM

Deleted: spread of the distributions

Rhys Whitley 11/5/16 2:54 PM

Deleted: highlight possible biases in the models (over- or underestimating the tower fluxes), as well as their ability to capture the spread of values

371 extremes (day-to-day variability) and the seasonality of the fluxes respectively,
372 highlighting possible predictive biases (i.e. the over- or underestimation of the tower
373 fluxes). Across the NATT, the PDFs for the tower fluxes tended to shift to low values and
374 became narrower as annual rainfall declined, and this was most prominent in the dry
375 season. A change in the spread and mean position of the flux tower PDFs demonstrate
376 the strong seasonality of water and carbon exchange at all sites. The PDFs of the model
377 simulations did not replicate this trend, having high densities and being mostly
378 stationary across sites. Regarding savanna water-use, the distributions of the BIOS2 and
379 SPA models were similar to those of the flux towers. The BESS model also showed a
380 similar distribution of LE, despite the fact that it did not simulate soil water extraction.
381 The LPJGUESS model, which had the shallowest simulated tree rooting depth, displayed
382 PDFs of high density that were biased towards low LE (20 – 40 W m⁻²) across all sites
383 and seasons. The MAESPA model showed a similar behaviour, despite this model having
384 a much deeper simulated rooting depth and a root-water extraction scheme that is
385 equivalent to the SPA model. The distributions for the CABLE and BIOS2 models were
386 largely disparate despite these models being functionally equivalent. Notably, CABLE
387 wet season LE was more broadly distributed (5 – 200 W m⁻²) than the flux towers and
388 other models at all sites, while dry season LE was narrower. In relation to savanna
389 carbon uptake, all models showed wet and dry season PDFs of high density that became
390 more closely aligned with the flux tower distributions as the sites became drier. The
391 behaviour of the modelled GPP distributions were otherwise similar to those of the
392 modelled LE distributions. The differences among TBM and flux tower PDFs indicated
393 possible issues in simulated processes that are active during the wet season.

394 The benchmarks set low to high levels of expected TBM performance across the NATT.
395 Additionally, they also demonstrated the level of model complexity that is required to
396 simulate water and carbon exchange at these sites. The simplest of the benchmarks,
397 represented as a linear regression of the fluxes against R_s (emp1), which was capable of
398 predicting the magnitude and daily time-course of the tower fluxes (data not shown),
399 but there was not enough information in R_s to capture the seasonality or the distribution
400 of the fluxes expressed by the tower data. The intermediate benchmark that included
401 additional meteorological information on T_a and VPD (emp2) demonstrated an
402 improved capability in capturing the flux distributions, but could not replicate the full
403 seasonality of the fluxes across the NATT. It was only by including additional
404 phenological information (LAI) together with site meteorology (R_s , T_a and VPD) that the
405 seasonality and distribution of the fluxes could be captured, as demonstrated by the

Rhys Whitley 11/5/16 3:00 PM

Deleted: T

Rhys Whitley 16/5/16 11:06 AM

Deleted: By contrast, t

Rhys Whitley 11/5/16 3:07 PM

Deleted: , especially for wet season GPP

410 most complex benchmark (emp3). This indicated that in order for the TBMs to achieve
411 the best possible performance at simulating water and carbon exchange along the NATT,
412 the correct implementation and utilisation of phenological information by the models
413 was required. All TBMs used in this study utilised this breadth of information, but only
414 some of the models were capable of meeting the expected level of performance set by
415 the emp3 benchmark, and only then for specific sites and seasons.

416

417 **3.2 Residual analysis**

418 An analysis of the model residuals was conducted to show how model structure affects
419 the prediction of savanna fluxes across the rainfall gradient. To do this we examined the
420 standardised model residuals from each TBM, determined by expressing the residual
421 error in terms of its standard deviation. Figure 4 shows the residual time-series for
422 model predicted LE and GPP at each savanna site and provides an effective way of
423 examining how a model responds to progressive changes in the environment, through
424 the expression of model bias and error (Medlyn et al., 2005).

425 The model residuals demonstrated that there was significant bias and heteroscedasticity
426 in predicted LE and GPP in almost all cases. The residual time-series showed that model
427 error was largest in the wet season, but declined with the transition into the dry season.
428 Additionally, the models underestimated LE and GPP more significantly during the wet
429 season . A possible explanation for this behaviour is that during the wet season, multiple
430 land-surface components: the soil surface, the understorey grasses, and the tree canopy
431 (i.e. 3 sources for potential error) contribute to the bulk fluxes, while during the dry
432 season only the tree canopy contributes (i.e. 1 source for potential error). It is likely that
433 the reduction in residual error between wet and dry seasons was a result of the
434 declining influence of the grasses and the soil surface to ecosystem land-surface
435 exchange during the latter period (via senescence and low surface soil moisture
436 respectively). The bias towards the underestimation of wet season fluxes was more
437 pronounced at the mesic sites (Howard Springs, Adelaide River), despite some models
438 simulating relatively deep root profiles (e.g. BIOS2, MAESPA). Differences in how the
439 TBMs simulated root-water extraction also had no effect on reducing this bias (e.g.
440 MAESPA, SPA). Given that soil-water was not a limiting factor at the mesic sites during
441 this period, deep root profiles offered limited advantage towards model performance.
442 Nonetheless, the simulated tree root-zone appeared to be an important factor for all

443 [sites during the dry season, with shallow root depths \(LPJGUESS: 2 m\) and/or](#)
444 [inadequate root-water uptake schemes \(CABLE: concentrated in the upper soil profile\)](#)
445 [the likely cause for underestimation during this period. However, as the sites became](#)
446 [drier \(e.g. Sturt Plains\) a shallow root-profile was suitable to give flux estimates of a](#)
447 [reasonably low error. Despite model error reducing with the increase in ecosystem](#)
448 [water limitation that occurs in both space \(down the NATT\) and time \(wet to dry](#)
449 [season\), there are still patterns of model bias that may be unrelated to simulated soil-](#)
450 [water dynamics. This is particularly obvious during the wet-to-dry transition periods](#)
451 [\(e.g. BIOS2, SPA\) when the C₄ grass understorey senesces, indicating possible problems](#)
452 [with the how the models translate information on phenology.](#)

453

454 **3.3 Model performance**

455 Figure 5 shows a comparison of individual TBM performance ordered by site from
456 wettest (Howard Springs) to driest (Sturt Plains) and in terms of their annual, wet and
457 dry season predictions for each flux. Despite differences in model complexity (Table 1),
458 the TBMs showed a similar performance across sites and seasons. For almost all sites,
459 the TBMs outperformed the emp1 benchmark for annual flux predictions (Fig. 5a).
460 However, there were some exceptions to this, and good performance in one flux did not
461 necessarily result in good performance in the other. For example, MAESPA was unable
462 to beat the emp1 benchmark for LE at sites where MAP > 1000 mm, but performed
463 better than the emp2 benchmark for GPP. In general, there was a slight pattern of
464 increased model performance as annual rainfall declined, though with a degree of site-
465 to-site variability in the rankings for some of the TBMs.

466 In order to examine how seasonal changes affect model performance, we additionally
467 determined the metrics and rankings for the wet and dry season periods (Fig. 5b-c).
468 Seasonal differences were immediately obvious. Model performance for wet season LE
469 and GPP was low to moderate, and the majority of the TBMs showed a performance that
470 ranged between the emp1 and emp2 benchmarks. In contrast, there were noticeable
471 improvements to dry season model performance amongst the TBMs. For dry season LE,
472 half the models (BIOS2, BESS, and SPA) were able to consistently outperform the emp2
473 benchmark, and come close to meeting the same number of metrics as the emp3
474 benchmark particularly at the drier sites. In comparison, predicted dry season GPP saw
475 a larger enhancement in model performance, with TBMS more frequently outperforming

476 the emp2 benchmark and even some outperforming the emp3 benchmark (LPJGUESS,
477 BESS, and SPA at the Daly Uncleared site). The exception to all this was the CABLE
478 model, which showed surprisingly little loss or gain in performance despite the season.
479 The results give an indication that as a whole, input information was better utilised by
480 each TBM at drier sites and in the dry season, suggesting that there are problems in wet
481 season processes.

482

483 **4. Discussion**

484 The NATT, which covers a marked rainfall gradient, presents a natural 'living laboratory'
485 with which a models ability to simulate fluxes in savanna ecosystems may be assessed.
486 Our results have highlighted that there is a clear failure of the models to adequately
487 perform at predicting wet season dynamics, as compared to the dry season, and
488 suggests that modelled processes relating to the C₄ grass understorey are insufficient.
489 This highlights a key weakness of this group of TBMs, which likely extends to other
490 models outside of this study. The inability of these TBMs to capture wet season
491 dynamics is highlighted by the benchmarking, where the performance for many of the
492 models was at best equivalent to that of a multi-linear regression against R_s , T_a and VPD
493 (emp2) and in some cases no better than a linear regression against R_s (emp1). Given
494 that this subset of TBMs are sophisticated process-based models that represent our best
495 understanding of land-surface, atmospheric exchange processes, we would expect them
496 to perform as well as a neural network prediction (emp3). Consequently there is an
497 evident underutilisation of the driving information (i.e. a failure to describe the
498 underlying relationships in the data) impeding the performance of these models when
499 predicting savanna fluxes. However, there were instances where some of the TBMs were
500 able to reach similar levels of performance with the emp3 benchmark, and strongly
501 suggests that each of these models is capable of replicating savanna dynamics under
502 certain conditions (e.g. during the dry season).

503

504 Our results suggest that errors among models are likely to be systematic, rather than
505 related to calibration of existing parameters. For example, BIOS2 had previously
506 optimised model parameters for Australian vegetation (see Haverd et al.2013), but was
507 still unable to out-perform the emp3 benchmark in most cases, although [it](#) performed
508 better than an un-calibrated CABLE, to which it is functionally similar. Similarly,
509 MAESPA and SPA, which used considerable site characteristic information to

510 parameterise their simulations, did not significantly outperform un-calibrated models
511 (e.g. CABLE). Additionally, despite these models using the same leaf, root and soil
512 parameterisations, both SPA and MAESPA displayed markedly different performances in
513 predicting LE. Consequently, improving how models represent key processes that drive
514 savanna dynamics is critical to improving model performance across this ecosystem.

515

516 There is certainly enough information in the time-varying model inputs to be able to
517 adequately simulate wet and dry season dynamics, as is evidenced by the benchmarks.

518 We therefore consider the implications of our results, and present possible reasons
519 below for why this group of TBMs is failing to capture water and carbon exchange along
520 the NATT, and make suggestions as to how this could be improved.

521

522 **4.1 Water access and tree rooting depth**

523 During the late dry season surface soil moisture in the sandy soils declines to less than
524 3% volumetric water content, with an equivalent matric potential of 3 to 4 MPa (Prior et
525 al., 1997). During this seasonal phase, the grass understorey becomes inactive and LE
526 can be considered as equivalent to tree transpiration, such that it is the only active
527 component during this period (O'Grady et al., 1999). Using this equivalence, one can
528 infer the relative effect that rooting depth has on LE during this period. Previous studies
529 have shown that for these savanna sites along the NATT, tree transpiration is
530 maintained throughout the dry season by deep root systems that access deep soil-water
531 stores, which in turn are recharged over the wet season (Eamus et al., 2000; Hutley et
532 al., 2001; Kelley et al., 2007; O'Grady et al., 1999). In order for models to perform well
533 they will need to set adequate rooting depths and distributions, along with root water
534 uptake process, to enable a model response to such seasonal variation. Examining
535 performance across the models, we can infer this to be a key deficiency. As expected,
536 TBMs that prescribed shallow rooting depths (e.g. LPJGUESS) did not simulate this
537 process well, and underestimated dry season LE at 3 of the 5 savanna sites by up to 30
538 to 40%. The two sites at Adelaide River and Sturt Plains were an exception to this with
539 the TBMs displaying a low residual error, which is likely to be a consequence of heavier
540 textured soils and trees at these sites having shallow root profiles. At Adelaide River
541 shallow root profiles are a consequence of shallow, heavier textured soils, however dry
542 season transpiration is sustained due to the presence of saturated yellow hydrosol soils.
543 | Sturt Plains is a grassland ([the](#) end member of the savanna continuum) where C₄ grasses
544 dominate and no trees are present such that transpiration is close to zero in the dry

Rhys Whitley 10/5/16 12:25 PM
Deleted: have

545 | season. The few small shrubs that **are** established have shallow root profiles that have
546 | adapted to isolated rainfall events driven by convective storms (Eamus et al., 2001;
547 | Hutley et al., 2001, 2011). Consequently, the TBMs would be expected to perform better
548 | at these sites, as water and carbon exchange will be modulated by the soil-water status
549 | of the sub-surface soil layers. For the other sites, models which assumed a root depth > 5
550 | m (BIOS2, SPA and MAESPA), showed the most consistent performance in predicting dry
551 | season LE, and we suggest for seasonally water-limited ecosystems, such as savanna,
552 | that deeper soil water access is critical. Our results highlight the need for data with
553 | which to derive more mechanistic approaches to setting rooting depth, such as that of
554 | Schymanski et al. (2009).

555 | Interestingly, a low residual error for LE in the dry season, did not translate as good
556 | performance in the overall model ranking. This suggests that other processes along the
557 | soil-vegetation-atmosphere continuum need to be considered to improve simulated
558 | woody transpiration. Such processes may include root-water uptake (distribution of
559 | roots and how water is extracted), and the effect of water stress and increased
560 | atmospheric demand at the leaf-level (adjustment of stomatal conductance due to
561 | changes in leaf water potential). More detailed model experiments that examine how
562 | each TBM simulates these processes would help identify how they can be improved.

563 | An exception to the above is the BESS model, which forgoes simulating belowground
564 | processes of soil hydrology and root-water uptake entirely. Rather, this model assumes
565 | that the effects of soil-moisture stress on water and carbon exchange is expressed
566 | through changes in LAI (and by extension V_{cmax}), which acts as a proxy for changes in soil
567 | moisture content (Ryu et al., 2011). The fact that BESS performed moderately well along
568 | the NATT, coupled with the fact that tree transpiration continues through the dry season
569 | suggests that there may be enough active green material for remote sensing proxies of
570 | water-stress to generally work rather well for savanna ecosystems. It is notable that
571 | BESS overestimated both GPP and ET in dry season at the driest site, Sturt Plains (Fig
572 | 2e), implying that greenness detected by satellite remote sensing might not capture
573 | carbon and water dynamics well in such a dry site.

574

575 | **4.2 Savanna wet season dynamics**

576 | The relative performance of the TBMs at predicting LE was much poorer in the wet
577 | season compared to the dry season. The reason for this difference is that wet season LE

579 is the sum of woody and herbaceous transpiration (E_{veg}) as well as soil and wet-surface
580 evaporation (E_{soil}); in contrast dry season LE is predominantly woody transpiration as
581 described previously. During the wet season, up to 75% of total LE arises from
582 understory herbaceous transpiration and soil evaporation (Eamus et al., 2001; Hutley
583 et al., 2000; Moore et al., 2016) and of this fraction the C_4 grasses contribute a significant
584 daily amount (Hutley et al., 2000). In the absence of observations of understory LE it can
585 be difficult to determine whether grass transpiration is being simulated correctly.
586 However, separating out the components of wet season LE into soil and vegetation can
587 help identify which of these components are causes for error.

588
589 Separating the outputs of simulated E_{veg} and E_{soil} from each TBM (excluding BESS which
590 did not determine these as outputs during the study) shows that simulated wet season
591 E_{veg} was particularly low for a lot of the models, despite high LAI and non-limiting soil-
592 water conditions (Figure 6). A previous study at Howard Springs by Hutley et al. (2000)
593 observed that during the wet season, the grass understory could transpire $\sim 2.8 \text{ mm d}^{-1}$,
594 while the tree canopy transpired only 0.9 mm d^{-1} ($E_{veg} = 3.7 \text{ mm d}^{-1}$). Of the 6 TBMs at
595 Howard Springs, only CABLE and SPA were able to predict an E_{veg} close to this level,
596 while the other models predicted values closer to tree transpiration (i.e. an under-
597 estimate). This pattern is similar for other NATT sites, where predicted wet season E_{veg}
598 remained low and was dominated by E_{soil} at the southern end of the NATT. An
599 underestimation of wet season LE could be due to underestimated E_{soil} in some of the
600 models. Conversely, CABLE and BIOS2 predicted a higher E_{soil} than the other models, and
601 this could be a reason for their higher LE performance during the wet season. Although
602 E_{soil} has been reported to reach as high as 2.8 mm d^{-1} at Howard Springs (Hutley et al.,
603 2000), predicted E_{soil} by these models may still be overestimated, given that vegetation
604 cover during this period is at a seasonal peak (limiting energy available at the soil
605 surface) and transpiration is only limited by available energy not water (Hutley et al.,
606 2000; Ma et al., 2013; Schymanski et al., 2009; Whitley et al., 2011). Given the limited
607 data for E_{soil} along the NATT, it is difficult to determine how large E_{soil} should be.
608 However, the ratios displayed by the TBMs appear to be reasonable though, with
609 vegetation acting as the predominant pathway for surface water flux.

610
611 Grass transpiration is thus clearly being under-represented by most of the TBMs, and
612 reasons for this could be due to multiple factors. The evolution of C_4 grasses to fix
613 carbon under low light, low CO_2 concentrations and high temperatures has resulted in a
614 gas-exchange process that is highly water-use efficient (von Caemmerer and Furbank,

Rhys Whitley 10/5/16 12:27 PM
Deleted: understorey

Rhys Whitley 10/5/16 12:27 PM
Deleted: s

Rhys Whitley 10/5/16 12:32 PM
Deleted: that we discuss below

1999). Consequently, this life form is abundant in tropical, water-limited ecosystems, where it can contribute to more than 50% of total LAI (2.0 to 2.5), particularly at high rainfall sites (Sea et al., 2011). The annual strategy of the C₄ grasses at these sites is to indiscriminately expend all available resources to maximise productivity during the monsoon period, for growth and to increase leaf area. This therefore allows grass transpiration to exceed tree transpiration during the peak wet season as evergreen trees will be more conservative in their water-use, allowing them to remain active in the dry season (Eamus et al., 2001; Hutley et al., 2000; Scholes and Archer, 1997). Following this logic, our results suggest that the TBMs are either: i) incorrectly ascribing leaf area to the understorey (i.e. the C₄ fractional cover is too low), ii) incorrectly describing the C₄ leaf-gas exchange physiology, iii) incorrectly describing the understorey micro climatic environment (R_s , T_a , VPD), or iv) a combination of these causes. Furthermore, it should be noted that the TBMs used in this study are not truly modelling grasses, but approximating them. Grasses are effectively simulated as 'stem-less' trees, and the distinction between the two life forms is reliant on different parameter sets (e.g. V_{cmax} , height, etc.) and slight modifications of the same process (e.g. rate of assimilation, respiration, etc.). While our results and the tower data do not allow us to directly determine how C₄ grasses may be misrepresented in these TBMs, they clearly indicate that future development and evaluation should be focused on these issues. Eddy covariance studies of understorey savanna vegetation as conducted by Moore et al. (2016) will be critical to this process.

4.3 Savanna phenology

The results from this study have shown that to simulate savanna fluxes, TBMs must be able to simulate the dynamics of savanna phenology, expressed by LAI. This was highlighted by the empirical benchmarks, where the results showed that while R_s , T_a and VPD were important drivers, LAI was required to capture the seasonality and magnitude of the fluxes to achieve good performance. LAI integrates the observed structural changes of the savanna as annual rainfall declines with reduced woody stem density, driving water and carbon exchange as a result (Kanniah et al., 2010; Ma et al., 2013; Sea et al., 2011). If LAI is prescribed in a model, it is important that leaf area is partitioned correctly between the trees and grass layers to describe their respective phenology. This partitioning is important, as the C₄ grass understorey explains most of the seasonal variation in LAI, and is a consequence of an annual phenology that exhibits rapid growth at the onset of the wet season and senescence at the onset of the dry (Williams et al.,

Rhys Whitley 10/5/16 12:33 PM
Deleted ;

654 1996b). By contrast the evergreen eucalypt canopy shows modest reductions in canopy
655 leaf area during the dry season, especially as mean annual rainfall declines (Bowman
656 and Prior, 2005; Kelley et al., 2007). The strong seasonal dynamics of the grasses result
657 in large changes in LAI, with levels varying between 0.7 and 2.5 at high rainfall sites (Sea
658 et al., 2011). The phenological strategy of the C₄ grasses also changes with rainfall
659 interannual variability, with the onset of the greening period becoming progressively
660 delayed as sites become drier, to become eventually rain-pulse driven as the monsoonal
661 influence weakens (Ma et al., 2013).

662

663 With the exception of LPJGUESS, all models prescribed LAI as an input driver.
664 Prescribing LAI can be problematic depending on the time-scale and how it is
665 partitioned between trees and grass layers. At large time-steps (months) it will fail to
666 capture the rapidly changing dynamics of vegetation during the transition periods, and
667 this is particularly true for the onset of the wet season (Sep-Nov) especially at drier sites
668 that are subject to larger interannual rainfall variability (Hutley et al., 2011).
669 Additionally, as the sites become drier the tree:grass ratio will become smaller and this
670 dynamic can be difficult to predict, although methods do exist (see Donohue et al. 2009).
671 From the results, we infer that TBMs that prescribe LAI and allow for a dynamic
672 representation of tree and grass ratios are better able to capture the changing dynamics
673 of the savanna system. This is a possible explanation for the better performance of the
674 BIOS2, MAESPA and SPA models in simulating GPP as these models dynamically
675 partition leaf area between trees and grasses at the sub-monthly time-scale, rather than
676 using a bulk value. However, there are limitations to using prescribed LAI,
677 predominantly in that it describes a stable system, of which savannas are typically not,
678 having a large sensitivity to changes in climate, particularly rainfall variability and
679 disturbance (Sankaran et al., 2005). DGVMs that consider dynamic vegetation and use a
680 prognostic LAI can simulate the feedback between the climate and the relative cover of
681 trees and grasses, which shapes the savanna continuum. This feedback allows the
682 simulated savanna structure to potentially shift to alternate states (e.g. grassland or
683 forest) in response to changes in annual rainfall and fire severity (Scheiter and Higgins,
684 2007, 2009). While LPJGUESS was the only TBM to use a prognostic LAI in our study, it
685 achieved only moderate performance, and this may be due to how carbon is allocated
686 from the pool on an annual time step, such that it is not as dynamic as it could be.
687 However, its capability to simulate the feedback between climate and LAI is critical for
688 simulating how savanna dynamics may change from year to year. There may also be
689 issues with how phenology is simulated, particularly as it is determined from empirical

Rhys Whitley 10/5/16 12:33 PM

Deleted: ;

Rhys Whitley 10/5/16 12:33 PM

Deleted: ;

692 formulations, which are: i) not specifically developed for savanna environments and ii)
693 calculated before the growing season begins. Such formulations are therefore not
694 mechanistic, and do not respond to actual season dynamics (e.g. limiting soil water), but
695 are empirically determined (Richardson et al., 2013).

696

697 **5. Conclusions**

698

699 This study set out to assess how well a set of functionally different, state-of-the-art
700 TBMs perform at predicting the bulk exchanges of carbon and water over savanna land
701 surfaces. Our model inter-comparison has identified key weaknesses in the assumptions
702 of biosphere-atmosphere processes, which do not hold for savanna environments. Our
703 benchmarking has identified low model performance by TBMs is likely a result of
704 incorrect assumptions related to: i) deep soil water access, ii) a systematic under-
705 estimation of the contribution of the grass understorey in the wet season, and iii) the
706 use of static phenology to represent dynamic vegetation. Our results showed that these
707 assumptions, as they currently exist in TBMs, are not wholly supported by 'observations'
708 of savanna water and carbon exchange and need to be addressed if more reliable
709 projections are to be made on how savannas respond to environmental change. Despite
710 this, our benchmarking has shown that all TBMs could potentially operate well for
711 savanna ecosystems, provided that the above issues are developed. We suggest that
712 further work investigates how particular processes in the models may be affecting
713 overall predicted water and carbon fluxes, and may include testing variable rooting
714 depths, alternate root-water uptake schemes and how these might affect leaf-level
715 outputs (e.g. stomatal conductance, leaf water potential) among TBMs, and different
716 phenology schemes. The issues highlighted here also have scope beyond savanna
717 environments, and are relevant to other water-limited ecosystems. The results from this
718 study provide a foundation for improving how savanna ecosystem dynamics are
719 simulated.

720

721 **Acknowledgements**

722

723 This study was conducted as part of the 'Australian Savanna Landscapes: Past, Present
724 and Future' project funded by the Australian Research Council (DP130101566). The
725 support, collection and utilization of data were provided by the OzFlux network
726 (www.ozflux.org.au) and Terrestrial Ecosystem Research Network (TERN)

Rhys Whitley 10/5/16 12:34 PM

Deleted: on

Rhys Whitley 10/5/16 12:37 PM

Formatted: Not Strikethrough

728 (www.tern.org.au), and funded by the ARC (DP0344744, DP0772981 and
729 DP130101566). PALS was partly funded by the TERN ecosystem Modelling and Scaling
730 infrAStructure (eMAST) facility under the National Collaborative Research
731 Infrastructure Strategy (NCRIS) 2013-2014 budget initiative of the Australian
732 Government Department of Industry. Rhys Whitley was supported through the ARC
733 Discovery Grant (DP130101566). Jason Beringer is funded under an ARC FT
734 (FT110100602). Vanessa Haverd's contribution was supported by the Australian
735 Climate Change Science Program. We acknowledge the support of the Australian
736 Research Council Centre of Excellence for Climate System Science (CE110001028).
737

738 **References:**

- 739 Abramowitz, G.: Towards a public, standardized, diagnostic benchmarking system for
740 land surface models, *Geosci. Model Dev.*, 5(3), 819–827, doi:10.5194/gmd-5-819-2012,
741 2012.
- 742 Abramowitz, G., Leuning, R., Clark, M. and Pitman, A.: Evaluating the Performance of
743 Land Surface Models, *J. Clim.*, 21(21), 5468–5481, doi:10.1175/2008JCLI2378.1, 2008.
- 744 Ball, J. T., Woodrow, I. E. and Berry, J. A.: A model predicting stomatal conductance and
745 its contribution to the control of photosynthesis under different environmental
746 conditions., in *Progress in Photosynthesis Research*, pp. 221–224, Martinus-Nijhoff
747 Publishers, Dordrecht, the Netherlands., 1987.
- 748 Bashtannyk, D. M. and Hyndman, R. J.: Bandwidth selection for kernel conditional
749 density estimation, *Comput. Stat. Data Anal.*, 36(3), 279–298, doi:10.1016/S0167-
750 9473(00)00046-3, 2001.
- 751 Beringer, J., Hutley, L. B., Tapper, N. J. and Cernusak, L. A.: Savanna fires and their impact
752 on net ecosystem productivity in North Australia, *Glob. Chang. Biol.*, 13(5), 990–1004,
753 doi:10.1111/j.1365-2486.2007.01334.x, 2007.
- 754 Beringer, J., Hutley, L. B., Abramson, D., Arndt, S. K., Briggs, P., Bristow, M., Canadell, J. G.,
755 Cernusak, L. A., Eamus, D., Evans, B. J., Fest, B., Goergen, K., Grover, S. P., Hacker, J.,
756 Haverd, V., Kanniah, K., Livesley, S. J., Lynch, A., Maier, S., Moore, C., Raupach, M., Russell-
757 Smith, J., Scheiter, S., Tapper, N. J. and Uotila, P.: Fire in Australian Savannas: from leaf to
758 landscape., *Glob. Chang. Biol.*, 11, 6641, doi:10.1111/gcb.12686, 2014.
- 759 Best, M. J., Abramowitz, G., Johnson, H. R., Pitman, A. J., Balsamo, G., Boone, A., Cuntz, M.,
760 Decharme, B., Dirmeyer, P. A., Dong, J., Ek, M., Guo, Z., Haverd, V., van den Hurk, B. J.,
761 Nearing, G. S., Pak, B., Peters-Lidard, C., Santanello Jr, J. A., Stevens, L. and Vuichard, N.:
762 The plumbing of land surface models: benchmarking model performance, *J.*
763 *Hydrometeorol.*, 16, 1425–1442, 2015.
- 764 Bond, W. J.: What Limits Trees in C4 Grasslands and Savannas?, *Annu. Rev. Ecol. Evol.*
765 *Syst.*, 39(1), 641–659, doi:10.1146/annurev.ecolsys.39.110707.173411, 2008.
- 766 Bowman, D. M. J. S. and Prior, L. D.: Why do evergreen trees dominate the Australian
767 seasonal tropics?, *Aust. J. Bot.*, 53(5), 379–399, doi:10.1071/BT05022, 2005.
- 768 von Caemmerer, S. and Furbank, R. T.: Modeling C4 Photosynthesis, in *C4 Plant Biology*,
769 edited by R. F. Sage and R. K. Monson, pp. 173–211, Academic Press, Toronto., 1999.

770 Cernusak, L. A., Hutley, L. B., Beringer, J., Holtum, J. A. M. and Turner, B. L.:
771 Photosynthetic physiology of eucalypts along a sub-continental rainfall gradient in
772 northern Australia, *Agric. For. Meteorol.*, 151(11), 1462–1470,
773 doi:10.1016/j.agrformet.2011.01.006, 2011.

774 Chapin, F. S., Woodwell, G. M., Randerson, J. T., Rastetter, E. B., Lovett, G. M., Baldocchi, D.
775 D., Clark, D. A., Harmon, M. E., Schimel, D. S., Valentini, R., Wirth, C., Aber, J. D., Cole, J. J.,
776 Goulden, M. L., Harden, J. W., Heimann, M., Howarth, R. W., Matson, P. A., McGuire, A. D.,
777 Melillo, J. M., Mooney, H. A., Neff, J. C., Houghton, R. A., Pace, M. L., Ryan, M. G., Running, S.
778 W., Sala, O. E., Schlesinger, W. H. and Schulze, E. D.: Reconciling carbon-cycle concepts,
779 terminology, and methods, *Ecosystems*, 9(7), 1041–1050, doi:10.1007/s10021-005-
780 0105-7, 2006.

781 Chen, X., Hutley, L. B. and Eamus, D.: Carbon balance of a tropical savanna of northern
782 Australia, *Oecologia*, 137(3), 405–16, doi:10.1007/s00442-003-1358-5, 2003.

783 Collatz, G. J., Ball, J. T., Grivet, C. and Berry, J. A.: Physiological and environmental
784 regulation of stomatal conductance, photosynthesis and transpiration: a model that
785 includes a laminar boundary layer, *Agric. For. Meteorol.*, 54, 107–136, 1991.

786 Collatz, G. J., Ribas-Carbo, M. and Berry, J. A.: Coupled photosynthesis-stomatal
787 conductance model for leaves of C4 plants, *Funct. Plant Biol.*, 19(5), 519 – 538, 1992.

788 Donohue, R. J., Mc Vicar, T. R. and Roderick, M. L.: Climate-related trends in Australian
789 vegetation cover as inferred from satellite observations, 1981-2006, *Glob. Chang. Biol.*,
790 15(4), 1025–1039, doi:10.1111/j.1365-2486.2008.01746.x, 2009.

791 Duursma, R. A. and Medlyn, B. E.: MAESPA: A model to study interactions between water
792 limitation, environmental drivers and vegetation function at tree and stand levels, with
793 an example application to [CO₂] ?? drought interactions, *Geosci. Model Dev.*, 5(4), 919–
794 940, doi:10.5194/gmd-5-919-2012, 2012.

795 Eamus, D., O'Grady, A. P. O. and Hutley, L. B.: Dry season conditions determine wet
796 season water use in the wet-tropical savannas of northern Australia, *Tree Physiol.*,
797 20(18), 1219–1226, 2000.

798 Eamus, D., Hutley, L. B. and O'Grady, A. P. O.: Daily and seasonal patterns of carbon and
799 water fluxes above a north Australian savanna, *Tree Physiol.*, 21(12-13), 977–88, 2001.

800 Eamus, D., Chen, X., Kelley, G. and Hutley, L. B.: Root biomass and root fractal analyses of
801 an open Eucalyptus forest in a savanna of north Australia, *Aust. J. Bot.*, 50, 31–41, 2002.

802 Eamus, D., Cleverly, J., Boulain, N., Grant, N., Faux, R. and Villalobos-Vega, R.: Carbon and
803 water fluxes in an arid-zone Acacia savanna woodland: An analyses of seasonal patterns
804 and responses to rainfall events, *Agric. For. Meteorol.*, 182-183, 225–238,
805 doi:10.1016/j.agrformet.2013.04.020, 2013.

806 Farquhar, G. D., von Caemmerer, S. and Berry, J. A.: A Biochemical Model of
807 Photosynthetic CO₂ Assimilation in Leaves of C₃ species, *Planta*, 149, 78–90, 1980.

808 Grace, J., Jose, J. S., Meir, P., Miranda, H. S. and Montes, R. A.: Productivity and carbon
809 fluxes of tropical savannas, *J. Biogeogr.*, 33(3), 387–400, doi:10.1111/j.1365-
810 2699.2005.01448.x, 2006.

811 Haverd, V., Raupach, M. R., Briggs, P. R., Canadell, J. G., Isaac, P., Pickett-Heaps, C.,
812 Roxburgh, S. H., Van Gorsel, E., Viscarra Rossel, R. A. and Wang, Z.: Multiple observation
813 types reduce uncertainty in Australia’s terrestrial carbon and water cycles,
814 *Biogeosciences*, 10(3), 2011–2040, doi:10.5194/bg-10-2011-2013, 2013.

815 Haverd, V., Smith, B., Raupach, M. R., Briggs, P. R., Nieradzik, L. P., Beringer, J., Hutley, L.
816 B., Trudinger, C. M. and Cleverly, J. R.: Coupling carbon allocation with leaf and root
817 phenology predicts tree-grass partitioning along a savanna rainfall gradient,
818 *Biogeosciences*, 12, 16313–16357, doi:10.5194/bgd-12-16313-2015, 2016.

819 Haxeltine, A. and Prentice, I. C.: A general model for the light-use efficiency of primary
820 production, *Funct. Ecol.*, 10(5), 551–561, doi:10.2307/2390165, 1996.

821 Higgins, S. I. and Scheiter, S.: Atmospheric CO₂ forces abrupt vegetation shifts locally,
822 but not globally, *Nature*, 488, 209–212, doi:10.1038/nature11238, 2012.

823 Hsu, K.: Self-organizing linear output map (SOLO): An artificial neural network suitable
824 for hydrologic modeling and analysis, *Water Resour. Res.*, 38(12), 1–17,
825 doi:10.1029/2001WR000795, 2002.

826 Huntingford, C. and Monteith, J. L.: The behaviour of a mixed-layer model of the
827 convective boundary layer coupled to a big leaf model of surface energy partitioning,
828 *Boundary-Layer Meteorol.*, 88(1), 87–101, doi:10.1023/A:1001110819090, 1998.

829 Hutchinson, M. F. and Xu, T.: ANUCLIM v6.1, Fenner School of Environment and Society,
830 Australian National University, Canberra, ACT., 2010.

831 Hutley, L. B., O’Grady, A. P. O. and Eamus, D.: Evapotranspiration from Eucalypt open-
832 forest savanna of Northern Australia, *Funct. Ecol.*, 14, 183–194, 2000.

833 Hutley, L. B., Grady, A. P. O. and Eamus, D.: Monsoonal influences on evapotranspiration

834 of savanna vegetation of northern Australia, *Oecologia*, 126(3), 434–443,
835 doi:10.1007/s004420000539, 2001.

836 Hutley, L. B., Beringer, J., Isaac, P. R., Hacker, J. M. and Cernusak, L. A.: A sub-continental
837 scale living laboratory: Spatial patterns of savanna vegetation over a rainfall gradient in
838 northern Australia, *Agric. For. Meteorol.*, 151(11), 1417–1428,
839 doi:10.1016/j.agrformet.2011.03.002, 2011.

840 Isbell, R.: *The Australian Soil Classification, Revised Ed.*, CSIRO Publishing, Collingwood,
841 Victoria., 2002.

842 Kanniah, K. D., Beringer, J. and Hutley, L. B.: The comparative role of key environmental
843 factors in determining savanna productivity and carbon fluxes: A review, with special
844 reference to northern Australia, *Prog. Phys. Geogr.*, 34(4), 459–490,
845 doi:10.1177/0309133310364933, 2010.

846 De Kauwe, M. G., Zhou, S.-X., Medlyn, B. E., Pitman, A. J., Wang, Y. P., Duursma, R. A. and
847 Prentice, I. C.: Do land surface models need to include differential plant species
848 responses to drought? Examining model predictions across a latitudinal gradient in
849 Europe, *Biogeosciences Discuss.*, 12(15), 12349–12393, doi:10.5194/bgd-12-12349-
850 2015, 2015.

851 Kelley, G., O’Grady, A. P., Hutley, L. B. and Eamus, D.: A comparison of tree water use in
852 two contiguous vegetation communities of the seasonally dry tropics of northern
853 Australia: The importance of site water budget to tree hydraulics, *Aust. J. Bot.*, 55(7),
854 700–708, doi:10.1071/BT07021, 2007.

855 Koch, G. W., Vitousek, P. M., Steffen, W. L. and Walker, B. H.: Terrestrial transects for
856 global change research, *Vegetation*, 121(1-2), 53–65, doi:10.1007/BF00044672, 1995.

857 Kowalczyk, E. A., Wang, Y. P. and Law, R. M.: The CSIRO Atmosphere Biosphere Land
858 Exchange (CABLE) model for use in climate models and as an offline model, Aspendale,
859 Victoria., 2006.

860 Lehmann, C. E. R., Prior, L. D. and Bowman, D. M. J. S.: Decadal dynamics of tree cover in
861 an Australian tropical savanna, , 601–612, doi:10.1111/j.1442-9993.2009.01964.x,
862 2009.

863 Lehmann, C. E. R., Anderson, T. M., Sankaran, M., Higgins, S. I., Archibald, S., Hoffmann, W.
864 A., Hanan, N. P., Williams, R. J., Fensham, R. J., Felfili, J., Hutley, L. B., Ratnam, J., San Jose,
865 J., Montes, R., Franklin, D., Russell-Smith, J., Ryan, C. M., Durigan, G., Hiernaux, P., Haidar,
866 R., Bowman, D. M. J. S. and Bond, W. J.: Savanna vegetation-fire-climate relationships

867 differ among continents., *Science* (80-.), 343, 548–552, doi:10.1126/science.1247355,
868 2014.

869 Leuning, R.: A critical appraisal of a combined stomatal-photosynthesis model for C3
870 plants, *Plant, Cell Environ.*, 18(4), 339–355, 1995.

871 Li, L., Wang, Y. P., Yu, Q., Pak, B., Eamus, D., Yan, J., Van Gorsel, E. and Baker, I. T.:
872 Improving the responses of the Australian community land surface model (CABLE) to
873 seasonal drought, *J. Geophys. Res. Biogeosciences*, 117(4), 1–16,
874 doi:10.1029/2012JG002038, 2012.

875 Ma, X., Huete, A., Yu, Q., Coupe, N. R., Davies, K., Broich, M., Ratana, P., Beringer, J., Hutley,
876 L. B., Cleverly, J., Boulain, N. and Eamus, D.: Spatial patterns and temporal dynamics in
877 savanna vegetation phenology across the North Australian Tropical Transect, *Remote
878 Sens. Environ.*, 139, 97–115, doi:10.1016/j.rse.2013.07.030, 2013.

879 McKenzie, N. N., Jacquier, D., Isbell, R. and Brown, K.: *Australian soils and landscapes : an
880 illustrated compendium*, CSIRO Publishing, Collingwood, Victoria., 2004.

881 Medlyn, B. E., Robinson, A. P., Clement, R. and McMurtrie, R. E.: On the validation of
882 models of forest CO2 exchange using eddy covariance data: some perils and pitfalls, *Tree
883 Physiol.*, 25, 839–857, 2005.

884 Medlyn, B. E., Duursma, R. A., Eamus, D., Ellsworth, D. S., Prentice, I. C., Barton, C. V. M.,
885 Crous, K. Y., De Angelis, P., Freeman, M. and Wingate, L.: Reconciling the optimal and
886 empirical approaches to modelling stomatal conductance, *Glob. Chang. Biol.*, 17(6),
887 2134–2144, doi:10.1111/j.1365-2486.2010.02375.x, 2011.

888 Moore, C. E., Beringer, J., Evans, B., Hutley, L. B., McHugh, I. and Tapper, N. J.: The
889 contribution of trees and grasses to productivity of an Australian tropical savanna,
890 *Biogeosciences*, 13, 2387–2403, doi:10.5194/bg-13-2387-2016, 2016.

891 O’Grady, A. P. O., Eamus, D. and Hutley, L. B.: Transpiration increases during the dry
892 season: patterns of tree water use in eucalypt open-forests of northern Australia., *Tree
893 Physiol.*, 19(9), 591–597, 1999.

894 Parton, W. J., Anderson, D. W., Cole, C. V and Stewart, J. W. B.: Simulation of soil organic
895 matter formation and mineralization in semiarid agroecosystems, in *Nutrient Cycling In
896 Agricultural Ecosystems*, vol. 23, pp. 533–550., 1983.

897 Pitman, A. J.: The evolution of, and revolution in, land surface schemes designed for
898 climate models, *Int. J. Climatol.*, 23(5), 479–510, doi:10.1002/joc.893, 2003.

899 Prior, L. D., Eamus, D. and Duff, G. A.: Seasonal Trends in Carbon Assimilation, Stomatal
900 Conductance, Pre-dawn Leaf Water Potential and Growth in *Terminalia ferdinandiana*, a
901 Deciduous Tree of Northern Australian Savannas, *Aust. J. Bot.*, 45, 53–69, 1997.

902 Richardson, A. D., Keenan, T. F., Migliavacca, M., Ryu, Y., Sonnentag, O. and Toomey, M.:
903 Climate change, phenology, and phenological control of vegetation feedbacks to the
904 climate system, *Agric. For. Meteorol.*, 169, 156–173,
905 doi:10.1016/j.agrformet.2012.09.012, 2013.

906 Russell-Smith, J. and Edwards, A. C.: Seasonality and fire severity in savanna landscapes
907 of monsoonal northern Australia, *Int. J. Wildl. Fire*, 15(4), 541–550,
908 doi:10.1071/WF05111, 2006.

909 Ryu, Y., Baldocchi, D. D., Kobayashi, H., van Ingen, C., Li, J., Black, T. A., Beringer, J., van
910 Gorsel, E., Knohl, A., Law, B. E. and Rouspard, O.: Integration of MODIS land and
911 atmosphere products with a coupled-process model to estimate gross primary
912 productivity and evapotranspiration from 1 km to global scales, *Global Biogeochem.*
913 *Cycles*, 25(GB4017), doi:10.1029/2011GB004053, 2011.

914 Ryu, Y., Baldocchi, D. D., Black, T. A., Detto, M., Law, B. E., Leuning, R., Miyata, A.,
915 Reichstein, M., Vargas, R., Ammann, C., Beringer, J., Flanagan, L. B., Gu, L., Hutley, L. B.,
916 Kim, J., McCaughey, H., Moors, E. J., Rambal, S. and Vesala, T.: On the temporal upscaling
917 of evapotranspiration from instantaneous remote sensing measurements to 8-day mean
918 daily-sums, *Agric. For. Meteorol.*, 152(1), 212–222,
919 doi:10.1016/j.agrformet.2011.09.010, 2012.

920 Sankaran, M., Hanan, N. P., Scholes, R. J., Ratnam, J., Augustine, D. J., Cade, B. S., Gignoux,
921 J., Higgins, S. I., Le Roux, X., Ludwig, F., Ardo, J., Banyikwa, F., Bronn, A., Bucini, G., Caylor,
922 K. K., Coughenour, M. B., Diouf, A., Ekaya, W., Feral, C. J., February, E. C., Frost, P. G. H.,
923 Hiernaux, P., Hrabar, H., Metzger, K. L., Prins, H. H. T., Ringrose, S., Sea, W., Tews, J.,
924 Worden, J. and Zambatis, N.: Determinants of woody cover in African savannas., *Nature*,
925 438, 846–849, doi:10.1038/nature04070, 2005.

926 Scheiter, S. and Higgins, S. I.: Partitioning of root and shoot competition and the stability
927 of savannas., *Am. Nat.*, 170(4), 587–601, doi:10.1086/521317, 2007.

928 Scheiter, S. and Higgins, S. I.: Impacts of climate change on the vegetation of Africa: an
929 adaptive dynamic vegetation modelling approach, *Glob. Chang. Biol.*, 15(9), 2224–2246,
930 doi:10.1111/j.1365-2486.2008.01838.x, 2009.

931 Scheiter, S., Higgins, S. I., Beringer, J. and Hutley, L. B.: Climate change and long-term fire

932 management impacts on Australian savanna, *Glob. Chang. Biol.*, Submitted , 1211–1226,
933 doi:10.1111/nph.13130, 2014.

934 Scholes, R. J. and Archer, S. R.: Tree-grass interactions in savannas, *Annu. Rev. Ecol. Syst.*,
935 28, 517–544, doi:10.1146/annurev.ecolsys.28.1.517, 1997.

936 Schymanski, S. J., Roderick, M. L., Sivapalan, M., Hutley, L. B. and Beringer, J.: A test of the
937 optimality approach to modelling canopy properties and CO₂ uptake by natural
938 vegetation., *Plant. Cell Environ.*, 30(12), 1586–98, doi:10.1111/j.1365-
939 3040.2007.01728.x, 2007.

940 Schymanski, S. J., Sivapalan, M., Roderick, M. L., Beringer, J. and Hutley, L. B.: An
941 optimality-based model of the coupled soil moisture and root dynamics, *Hydrol. Earth
942 Syst. Sci.*, 12(3), 913–932, doi:10.5194/hess-12-913-2008, 2008.

943 Schymanski, S. J., Sivapalan, M., Roderick, M. L., Hutley, L. B. and Beringer, J.: An
944 optimality-based model of the dynamic feedbacks between natural vegetation and the
945 water balance, *Water Resour. Res.*, 45(1), doi:10.1029/2008WR006841, 2009.

946 Sea, W. B., Choler, P., Beringer, J., Weinmann, R. a., Hutley, L. B. and Leuning, R.:
947 Documenting improvement in leaf area index estimates from MODIS using
948 hemispherical photos for Australian savannas, *Agric. For. Meteorol.*, 151(11), 1453–
949 1461, doi:10.1016/j.agrformet.2010.12.006, 2011.

950 Simioni, G., Roux, X. Le, Gignoux, J. and Sinoquet, H.: Treegrass: a 3D, process-based
951 model for simulating plant interactions in tree–grass ecosystems, *Ecol. Modell.*, 131, 47–
952 63, 2000.

953 Simioni, G., Gignoux, J. and Le Roux, X.: Tree layer spatial structure can affect savanna
954 production and water budget: Results of a 3-D model, *Ecology*, 84(7), 1879–1894,
955 doi:10.1890/0012-9658(2003)084[1879:TLSSCA]2.0.CO;2, 2003.

956 Sitch, S., Smith, B., Prentice, I. C., Arneeth, A., Bondeau, A., Cramer, W., Kaplan, J., Levis, S.,
957 Lucht, W., Sykes, M., Thonicke, K. and Venevski, S.: Evaluation of ecosystem dynamics,
958 plant geography and terrestrial carbon cycling in the LPJ dynamic vegetation model,
959 *Glob. Chang. Biol.*, 9, 161–185, 2003.

960 Smith, B., Prentice, I. C. and Sykes, M. T.: Representation of vegetation dynamics in the
961 modelling of terrestrial ecosystems: Comparing two contrasting approaches within
962 European climate space, *Glob. Ecol. Biogeogr.*, 10(6), 621–637, doi:10.1046/j.1466-
963 822X.2001.t01-1-00256.x, 2001.

964 Still, C. J., Berry, J. a., Collatz, G. J. and DeFries, R. S.: Global distribution of C3 and C4
965 vegetation: Carbon cycle implications, *Global Biogeochem. Cycles*, 17(1), 6–1–6–14,
966 doi:10.1029/2001GB001807, 2003.

967 Wang, Y. P., Kowalczyk, E., Leuning, R., Abramowitz, G., Raupach, M. R., Pak, B., Van
968 Gorsel, E. and Luhar, A.: Diagnosing errors in a land surface model (CABLE) in the time
969 and frequency domains, *J. Geophys. Res. Biogeosciences*, 116, 1–18,
970 doi:10.1029/2010JG001385, 2011.

971 van der Werf, G. R., Randerson, J. T., Giglio, L., Gobron, N. and Dolman, A. J.: Climate
972 controls on the variability of fires in the tropics and subtropics, *Global Biogeochem.*
973 *Cycles*, 22(3), 1–13, doi:10.1029/2007GB003122, 2008.

974 Whitley, R. J., Macinnis-Ng, C. M. O., Hutley, L. B., Beringer, J., Zeppel, M., Williams, M.,
975 Taylor, D. and Eamus, D.: Is productivity of mesic savannas light limited or water
976 limited? Results of a simulation study, *Glob. Chang. Biol.*, 17(10), 3130–3149,
977 doi:10.1111/j.1365-2486.2011.02425.x, 2011.

978 Williams, M., Rastetter, E. B., Fernandes, D. N., Goulden, M. L., Wofsy, S. C., Shaver, G. R.,
979 Melillo, J. M., Munger, J. W., Fan, S.-M. and Nadelhoffer, K. J.: Modelling the soil-plant-
980 atmosphere continuum in a *Quercus Acer* stand at Harvard Forest: the regulation of
981 stomatal conductance by light, nitrogen and soil/plant hydraulic properties, *Plant, Cell*
982 *Environ.*, 19, 911–927, 1996a.

983 Williams, R. J., Duff, G. A., Bowman, D. M. J. S. and Cook, G. D.: Variation in the
984 composition and structure of tropical savannas as a function of rainfall and soil texture
985 along a large-scale climatic gradient in the Northern Territory, Australia, *J. Biogeogr.*,
986 23(6), 747–756, doi:10.1111/j.1365-2699.1996.tb00036.x, 1996b.

987

	Howard Springs ^a	Adelaide River ^b	Daly Uncleared ^c	Dry River ^d	Sturt Plains ^e
Years (inclusive)	2001 – 2012	2007 – 2009	2008 – 2012	2008 – 2012	2008 – 2012
Co-ordinates	12°29'39.12" S 131°09'09" E	13°04'36.84" S 131°07'04.08" E	14°09'33.12" S 131°23'17.16" E	15°15'31.62" S 132°22'14.04" E	17°09'02.76" S 133°21'01.14" E
Elevation (m)	64	90	110	175	250
^fMeteorology					
Annual Rainfall (mm)	1714	1460	1170	850	535
Min/Max Daily Temperature (°C)	22.0/33.0	21.8/35.3	20.8/35.0	20.0/34.8	19.0/34.2
Min/Max Absolute Humidity (g m ⁻³)	11.0/18.5	8.9/17.7	8.6/15.1	7.8/12.3	6.1/9.0
Min/Max Soil Moisture (m ³ m ⁻³)	0.06/0.1	0.09/0.14	0.03/0.06	0.03/0.05	0.04/0.1
Soil Temperature (°C)	32.7	35.7	32.8	<i>n.a.</i>	30.2
Solar Radiation (W m ⁻²)	256.5	258.1	270.6	266.5	269.7
Bowen Ratio	1.7	3.1	3.2	4.6	15.8
^fVegetation					
Overstorey species	<i>Eu. Miniata</i> <i>Eu. tetradonta</i> <i>Er. chlorostachys</i>	<i>Eu. tectifera</i> <i>Pl. careya</i> <i>Co. latifolia</i>	<i>Te. grandiflora</i> <i>Eu. tetradonta</i> <i>Co. latifolia</i>	<i>Eu. tetradonta</i> <i>Co. terminalis</i> <i>Eu. dichromophloia</i>	<i>n.a.</i>
Understorey species	<i>Sorghum</i> spp. <i>He. triticeus</i>	<i>Sorghum</i> spp. <i>Ch. fallax</i>	<i>Sorghum</i> spp. <i>He. triticeus</i>	<i>Sorghum intrans</i> <i>Th. Tiandra</i> <i>Ch. fallax</i>	<i>Astrabla</i> spp.
Basal Area (m ² ha ⁻¹)	9.7	5.1	8.3	5.4	<i>n.a.</i>
Canopy Height (m)	18.9	12.5	16.4	12.3	0.2
LAI (m ² m ⁻²)	1.04 ± 0.07	0.68 ± 0.07	0.80 ± 0.12	0.58 ± 0.11	0.39 ± 0.11
Total Leaf Nitrogen (g m ⁻³)	1.42 ± 0.20	1.27 ± 0.18	1.35 ± 0.19	1.97 ± 0.15	2.37 ± 0.17
^gSoil					
Type	Red kandosol	Yellow hydrosol	Red kandosol	Red kandosol	Grey vertosol
A Horizon	Texture Sandy loam	Sandy loam	Loam	Clay	loam
Clay PSD (%)	15	20	20	50	20
Sand PSD (%)	60	50	40	25	40
Thickness (m)	0.30	0.30	0.20	0.15	0.20
Bulk Density (Mg m ⁻³)	1.29	1.60	1.39	1.20	1.39
Hydraulic Conductivity (mm hr ⁻¹)	9	7	9	3	9
Field Capacity (mm m ⁻¹)	156	132	147	140	147
B Horizon	Texture Clay loam	Clay	Clay loam	Clay	Clay loam
Clay PSD (%)	40	55	35	55	35
Sand PSD (%)	30	20	30	20	30
Thickness (m)	1.20	0.60	0.69	1.29	0.69
Bulk Density (Mg m ⁻³)	1.39	1.70	1.39	1.39	1.39
Hydraulic Conductivity (mm hr ⁻¹)	8	5	7	2	7
Field Capacity (mm m ⁻¹)	146	31	146	107	146

988
989 **Table 1:** Summarised dataset information for each of the five savanna sites used in this study. This includes site descriptions pertaining to
990 local meteorology, vegetation and below ground soil characteristics. Where data were not available, the abbreviation *n.a.* is used. Definitions
991 for the species genus mentioned in the table are as follows: *Eucalyptus* (*Eu.*), *Erythrophleum* (*Er.*), *Terminalia* (*Te.*), *Corymbia* (*Co.*),
992 *Planchonia* (*Pl.*), *Buchanania* (*Bu.*), *Themda* (*Th.*), *Hetropogan* (*He.*), and *Chrysopogon* (*Ch.*). Eddy covariance datasets relating to each of the
993 5 sites here can be download from www.ozflux.org.au and hdl references are given by order of column (Jason Beringer (2013) – ^ahdl:
994 102.100.100/14228, ^bhdl: 102.100.100/14239, ^chdl: 102.100.100/14229, ^dhdl: 102.100.100/14234, ^ehdl: 102.100.100/14230). Site
995 meteorology is given as 30 year averages with values taken from ^fHutley, et al. (2011). Soil descriptions are taken from the Digital Atlas of
996 Australian Soils (www.asris.csiro.au) ^gIsbell, (2002).

Model Name	SPA	MAESPA	CABLE	BIOS2	BESS	LPJGUESS
Model definition	Soil-Plant-Atmosphere Model	MAESTRA-SPA	Community Atmosphere Biosphere Land-surface Exchange Model	Modified CABLE (CABLE + SLI + CASA-CNP)	Breathing Earth System Simulator	Lund-Potsdam-Jena General Ecosystem Simulator
Version	1.0	1.0	2.0	2.0	1.0	2.1
Reference	Williams et al. (1996a)	Duursma & Medlyn (2012)	Kowalyzck et al. (2006), Wang et al. (2011)	Haverd et al. (2013)	Ryu et al. (2011, 2012)	Smith et al. (2001)
Temporal resolution	30-min	30-min	30-min	Daily (30-min time-steps are generated from daily time-series)	Snap shot with MODIS overpass, then up-scaled to a daily and 8-day time series	Daily
Spatial resolution	Point	Point	0.05° (5 km)	0.05° (5 km)	0.05° (5 km)	Patch (c. 0.1 ha)
Functional class	Stand model	Individual Plant or Stand Model	Land-Surface Model	Land-Surface Model	Remote Sensing Model	Dynamic Global Vegetation Model
Canopy Description						
C₃ Assimilation	Farquhar et al. (1980)	Farquhar et al. (1980)	Farquhar et al. (1980)	Farquhar et al. (1980)	Farquhar et al. (1980)	Collatz et al. (1991)
C₄ Assimilation	Collatz et al. (1992)	Collatz et al. (1992)	Collatz et al. (1992)	Collatz et al. (1992)	Collatz et al. (1992)	Collatz et al. (1992)
Stomatal conductance	Williams et al. (1996a)	Medlyn et al. (2011)	Leuning (1995)	Leuning (1995)	Ball et al. (1987)	Haxeltine & Prentice (1996)
Transpiration	Penman-Monteith calculated at leaf-scale accounting for g_b and limitation of soil-water supply via ψ_l	Penman-Monteith calculated at the leaf scale	Penman-Monteith	Penman-Monteith	Penman-Monteith	Haxeltine & Prentice (1996)
Boundary layer resistance	$f(\text{wind speed, leaf width, air temperature})$	$f(\text{wind speed, leaf width, air temperature and atmospheric pressure})$	$f(\text{wind speed, leaf width, air temperature})$	$f(\text{wind speed, leaf width, air temperature})$	Not Modelled	Huntingford & Monteith (1998)
Aerodynamic resistance	$f(\text{wind speed, canopy height})$	Not calculated unless transpiration is calculated at the canopy scale, in which case g_b above isn't calculated.	$f(\text{wind speed, canopy height})$	$f(\text{wind speed, canopy height})$	$f(\text{wind speed, canopy height})$	Huntingford & Monteith (1998)
Leaf area index	Prescribed (MODIS)	Prescribed (MODIS)	Prescribed (MODIS)	Prescribed (MODIS)	Prescribed (MODIS)	Prognostic (C allocation)
Canopy structure	Canopy + understorey divided into 10 layers	Individual plant crowns, spatially explicit locations and uniform understorey	2 (tree/grass) big leaf (sunlit/shaded)	2 (tree/grass) big leaf (sunlit/shaded)	2 (tree/grass) big leaf (sunlit/shaded)	5-year age/size cohorts for trees, single-layer grass understorey
C₃:C₄ fraction	Dynamic ratio variable with time. Compete for water and light.	Dynamic ratio variable with time. Compete for water and light.	Simulated as independent layers	Dynamic ratio variable with time. Compete for water not light.	Still et al. (2003) Ratio changes 70:30 to 10:90 down transect	Prognostic, determined as the outcome of the competition with trees
Canopy interception	YES	YES	YES	YES	NO	YES
Simulates growth	NO	NO	NO	NO	NO	YES
Soil Profile Description						
Soil profile structure	Profile divided into N layers (prescribed - 20 in this case.)	Profile divided into N layers (prescribed - 20 in this case.)	Profile divided into 6 layers	Profile divided into 12 layers (adjustable)	Not Modelled	2 layers (0-0.5, 0.5-2 m) with 10 cm evaporation sub-layer
Soil hydraulic properties	Function of sand and clay particle size distributions	Function of sand and clay particle size distributions	Prescribed	Australian Soils Resource Information System (ASRIS)	Not Modelled	Sitch et al. (2003)
Soil depth	6.5 m	5.0 m	4.5 m	10.0 m	Not Modelled	2 m
Root depth	6.5 m	5.0 m	4.5 m	0.5 m (grasses), 5.0 m (trees)	Not Modelled	2 m
Root distribution	Prescribed; exponential decay as a function of surface biomass and the total root biomass of the column	Prescribed; exponential decay as a function of surface biomass and the total root biomass of the column	Prescribed; exponential decay	Prescribed; exponential decay	Not Modelled	PFT-specific, trees have deeper roots on average
Soil-water stress modifier	E_t via g_s is increased to meet atmospheric demand while ψ_l remains above a critical threshold	Maximum transpiration rate calculated from hydraulic conductance (soil-to-leaf) sets limit on actual transpiration, OR uses the Tuzet et al. (2003) model of stomatal conductance	Supply/Demand	g_s scaled by a soil moisture limitation function related to extractible water accessible by roots	Assumes LAI and seasonal variation of V_{cmax} reflect soil water stress	Supply/Demand
Hydraulic pathway resistance	$R_{soil} + R_{plant}$	$R_{soil} + R_{plant}$	Not Modelled	Not Modelled	Not Modelled	Not explicit, min(supply, demand) determines sapflow

997
998 **Table 2:** Summary table of the ecosystem models used in the experiment; highlighting differences and similarities in model structure and
999 shared processes. Information is broken down into how each model describes aboveground canopy and belowground soil processes.
1000

Statistical Metric	Definition
Correlation coefficient (r)	$\frac{n \sum_{i=1}^n (O_i M_i) - \sum_{i=1}^n O_i \sum_{i=1}^n M_i}{\sqrt{\left(n \sum_{i=1}^n O_i^2 - \left(\sum_{i=1}^n O_i \right)^2 \right) \left(n \sum_{i=1}^n M_i^2 - \left(\sum_{i=1}^n M_i \right)^2 \right)}}$
Standard Deviation (sd)	$\left 1 - \frac{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (M_i - \bar{M})^2}}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (O_i - \bar{O})^2}} \right $
Normalised mean error (NME)	$\frac{\sum_{i=1}^n M_i - O_i }{\sum_{i=1}^n \bar{O} - O_i }$
Normalised mean bias (MBE)	$\frac{1}{n} \sum_{i=1}^n (M_i - O_i)$

1001

1002 **Table A1:** Definition of common metrics used to determine ranks against the empirical benchmarks.

1003 The terms M and O stand for model and observations respectively, while n denotes the length of the
 1004 data, and i is the datum.

1005

1006 **Figure Captions**

1007

1008 **Figure 1:** The Northern Territory of Australia and the North Australian Tropical Transect (NATT)
1009 showing (a) the flux site locations with an accompanying 30-year (1970 to 2000) expression of the
1010 average meteorological conditions for (b) mean annual temperature, and (c) total annual
1011 precipitation derived from ANUCLIM v6.1 climate surfaces (Hutchinson and Xu, 2010).

1012

1013 **Figure 2:** Time-series of daily mean latent heat (LE) flux and gross primary productivity (GPP)
1014 depicting an average year for each of the 5 savanna sites using a smoothed, 7-day moving average.
1015 The sites are ordered from wettest to driest; (a) Howard Springs, (b) Adelaide River, (c) Daly River,
1016 (d) Dry River and (e) Sturt Plains. The joined, black dots are the tower flux time-series, while the
1017 grey lines are the performance benchmarks (emp1, emp2, emp3). Predictions of LE and GPP for each
1018 of the six terrestrial biosphere models are given by a spectrum of colours described in the legend.

1019

1020 **Figure 3:** Probability densities (expressed in scientific notation) of daily mean latent heat (LE) flux
1021 and gross primary productivity (GPP) at each of the 5 savanna sites, where the distributions for each
1022 flux are partitioned into wet and dry seasons. The order of the sites are from wettest to driest; (a)
1023 Howard Springs, (b) Adelaide River, (c) Daly River, (d) Dry River and (e) Sturt Plains. The grey
1024 region is the tower flux, while the dotted lines are the empirical benchmarks. Predicted LE and GPP
1025 probability densities from each of the six process-based models are given by a spectrum of colours
1026 described in the legend.

1027

1028 **Figure 4:** Standardised model residuals for latent energy (LE) and gross primary productivity (GPP)
1029 expressed in units of standard deviations (sd) $[(\text{modelled flux} - \text{observed flux})/\text{sd}(\text{observed flux})]$.
1030 Residuals are presented for each model: (a) CABLE, (b) BIOS2, (c) LPJGUESS, (d) MAESPA, (d) BESS
1031 and (e) SPA, where each flux site is represented by a blue-green-yellow gradient. For both fluxes, the
1032 residuals are plotted against time (ensemble average year) and against the flux prediction (bias).

1033

1034 **Figure 5:** Average rank plot showing the performance of the terrestrial biosphere models for all
1035 sites across the North Australian Tropical Transect (NATT) ordered in terms of annual rainfall as
1036 follows: Howard Springs (HowSpr), Adelaide River (AdrRiv), Daly Uncleared (DalUnc), Dry River
1037 (DryRiv), and Sturt Plains (StuPla). Models are individually ranked against the benchmarks in order
1038 of 1 to 4 (1 model + 3 benchmarks) and express the amount of metrics the models are meeting listed
1039 in Table S1. The rankings are determined individually for latent energy (LE) and gross primary
1040 productivity (GPP). The coloured lines represent each of the 6 models in the study, while the grey

1041 lines represent the empirical benchmarks. The average ranking for each model was determined for
1042 (a) a complete year, (b) the wet season and (c) the dry season.

1043

1044 **Figure 6:** Average year outputs of vegetation transpiration (grass + trees) and soil evaporation, as
1045 well as their percentage contributions to total latent energy (LE) for each of the 6 terrestrial
1046 biosphere models at each of the 5 savanna sites.

1047

1048 **Figure S1:** A smoothed (7-day moving average) representation of the environmental drivers used to
1049 construct the empirical benchmarks at each of the 5 savanna sites, and are shown from wettest to
1050 driest; (a) Howard Springs, (b) Adelaide River, (c) Daly River, (d) Dry River and (e) Sturt Plains. The
1051 time-series represents the seasonality over an average year for mean daily solar radiation (R_s), mean
1052 daily air temperature (T_a), mean daily vapour pressure deficit (VPD) and leaf area index (LAI).

1053

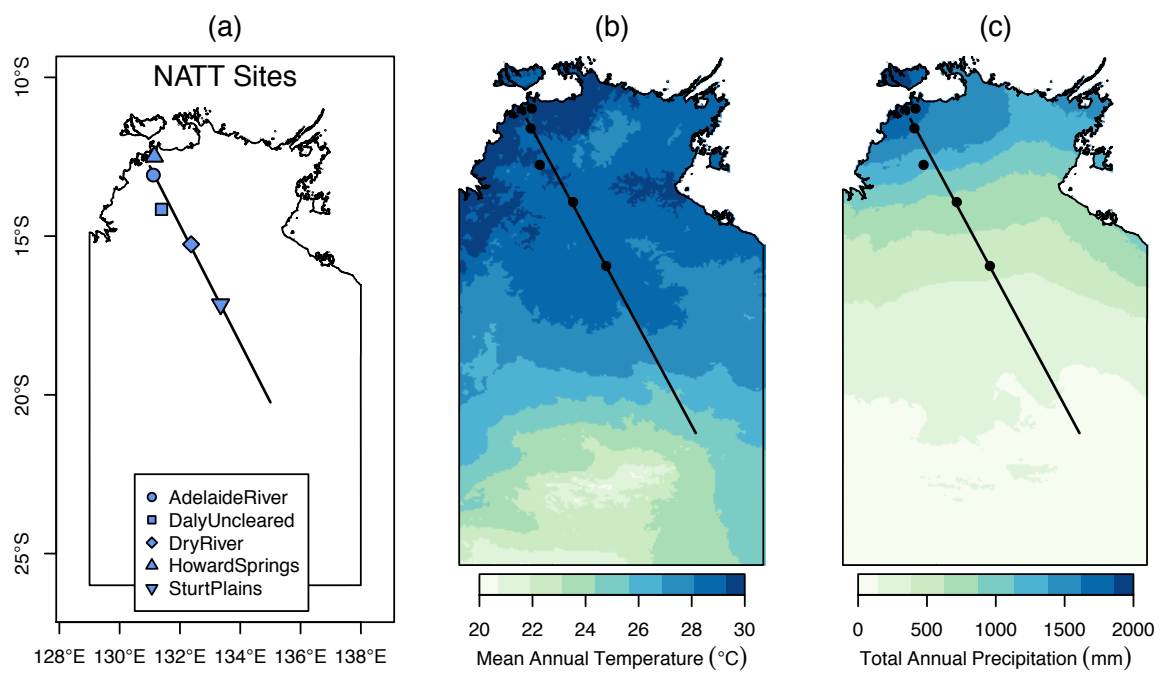


Figure 1: The Northern Territory of Australia and the North Australian Tropical Transect (NATT) showing (a) the flux site locations with an accompanying 30-year (1970 to 2000) expression of the average meteorological conditions for (b) mean annual temperature, and (c) total annual precipitation derived from ANUCLIM v6.1 climate surfaces (Hutchinson and Xu, 2010).

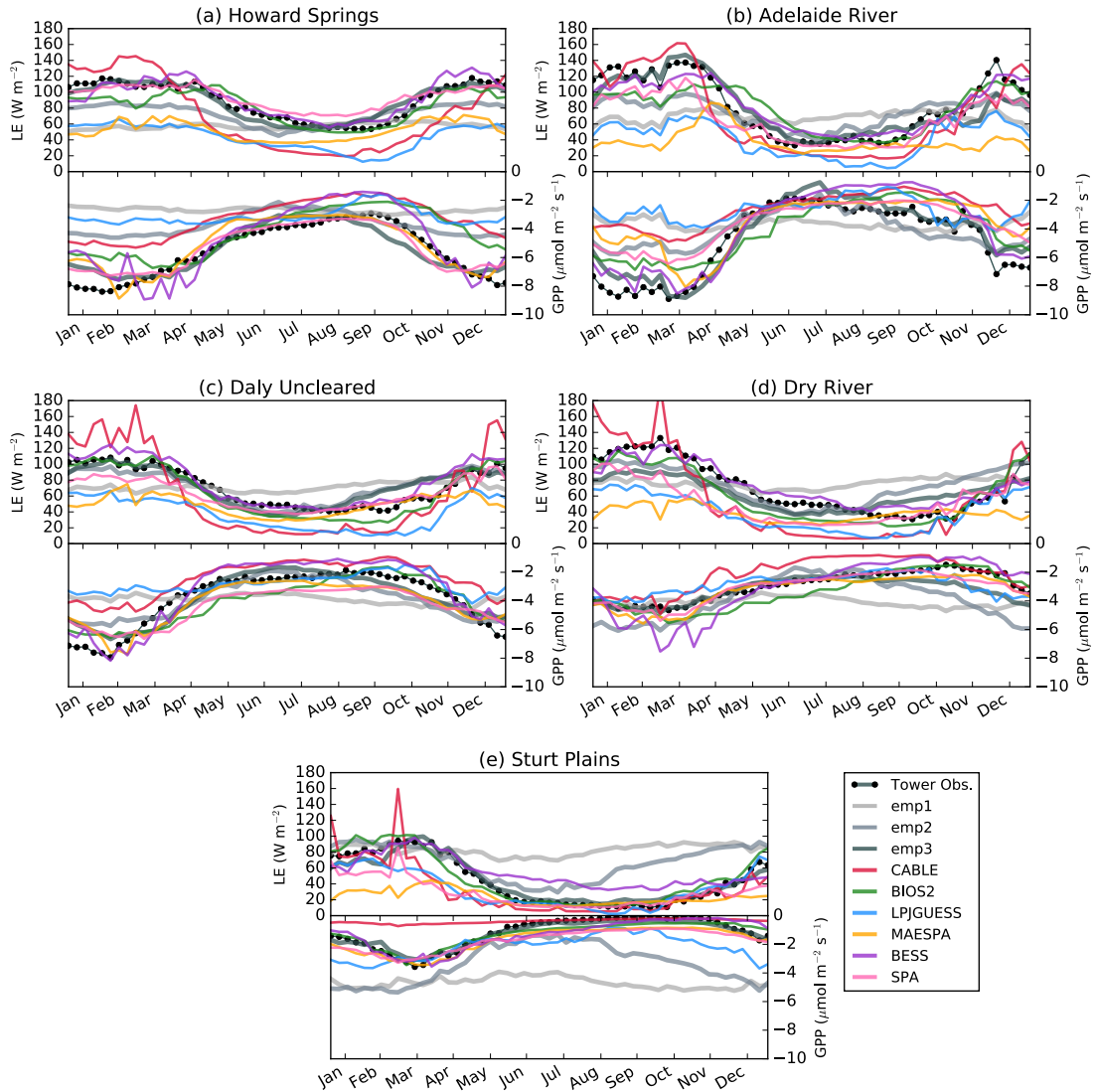


Figure 2: Time-series of daily mean latent heat (LE) flux and gross primary productivity (GPP) depicting an average year for each of the 5 savanna sites using a smoothed, 7-day moving average. The sites are ordered from wettest to driest; (a) Howard Springs, (b) Adelaide River, (c) Daly River, (d) Dry River and (e) Sturt Plains. The joined, black dots are the tower flux time-series, while the grey lines are the performance benchmarks (emp1, emp2, emp3). Predictions of LE and GPP for each of the six terrestrial biosphere models are given by a spectrum of colours described in the legend.

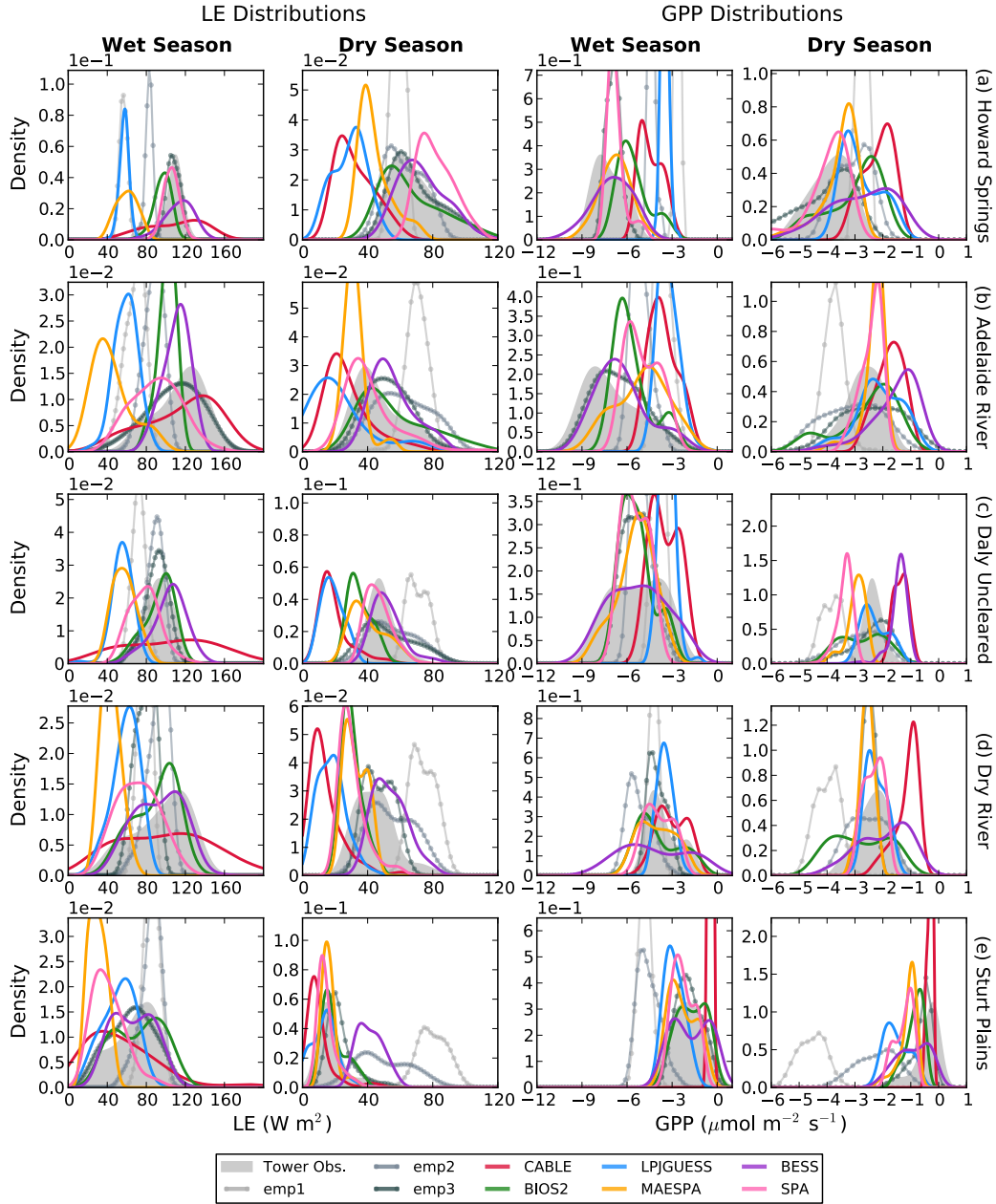


Figure 3: Probability densities (expressed in scientific notation) of daily mean latent heat (LE) flux and gross primary productivity (GPP) at each of the 5 savanna sites, where the distributions for each flux are partitioned into wet and dry seasons. The order of the sites are from wettest to driest; (a) Howard Springs, (b) Adelaide River, (c) Daly River, (d) Dry River and (e) Sturt Plains. The grey region is the tower flux, while the dotted lines are the empirical benchmarks. Predicted LE and GPP probability densities from each of the six process-based models are given by a spectrum of colours described in the legend.

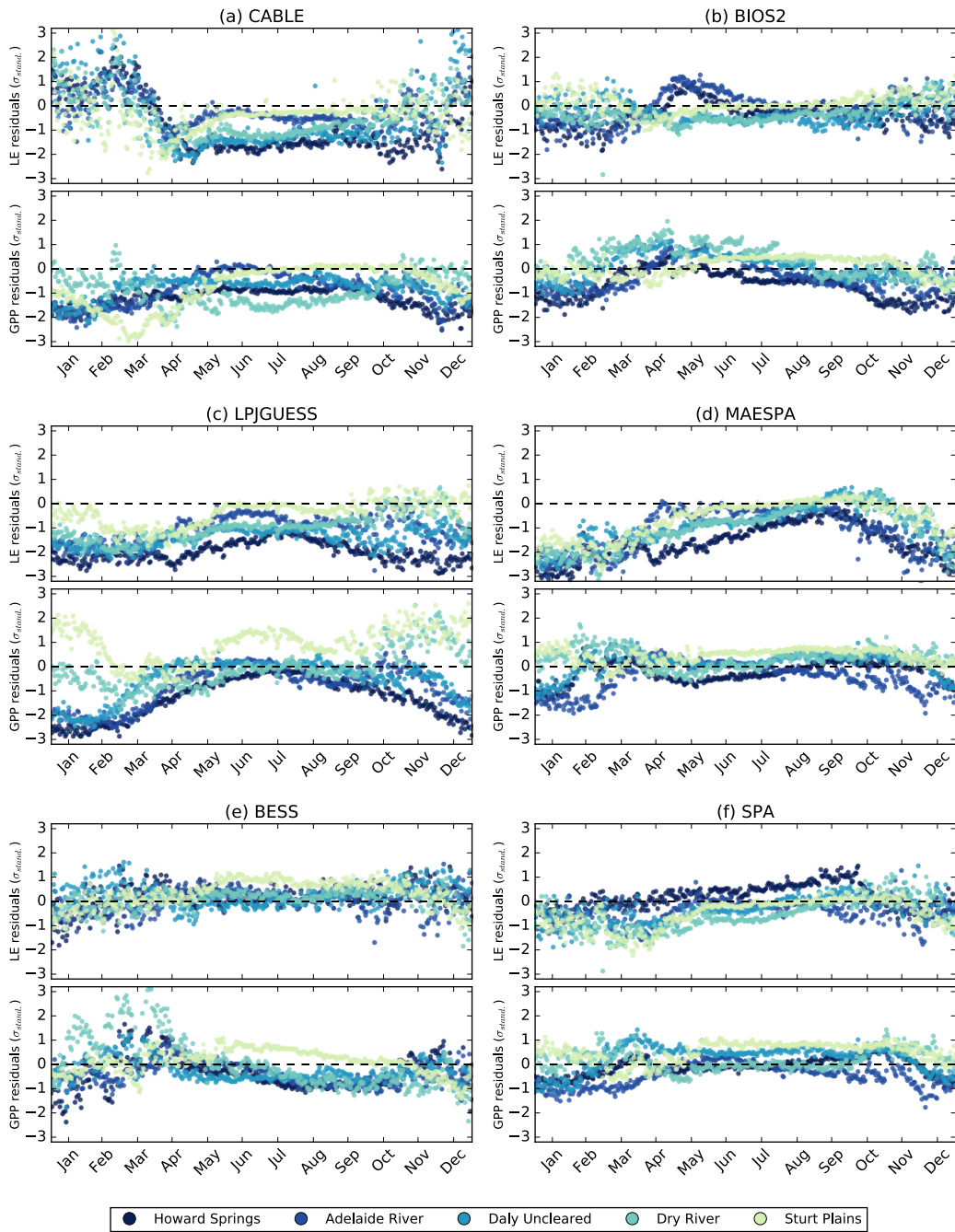


Figure 4: Standardised model residuals for latent energy (LE) and gross primary productivity (GPP) expressed in units of standard deviations (sd) $[(\text{modelled flux} - \text{observed flux})/\text{sd}(\text{observed flux})]$. Residuals are presented for each model: (a) CABLE, (b) BIOS2, (c) LPJGUESS, (d) MAESPA, (d) BESS and (e) SPA, where each flux site is represented by a blue-green-yellow gradient. For both fluxes, the residuals are plotted against time (an average year) and against the flux prediction (bias).

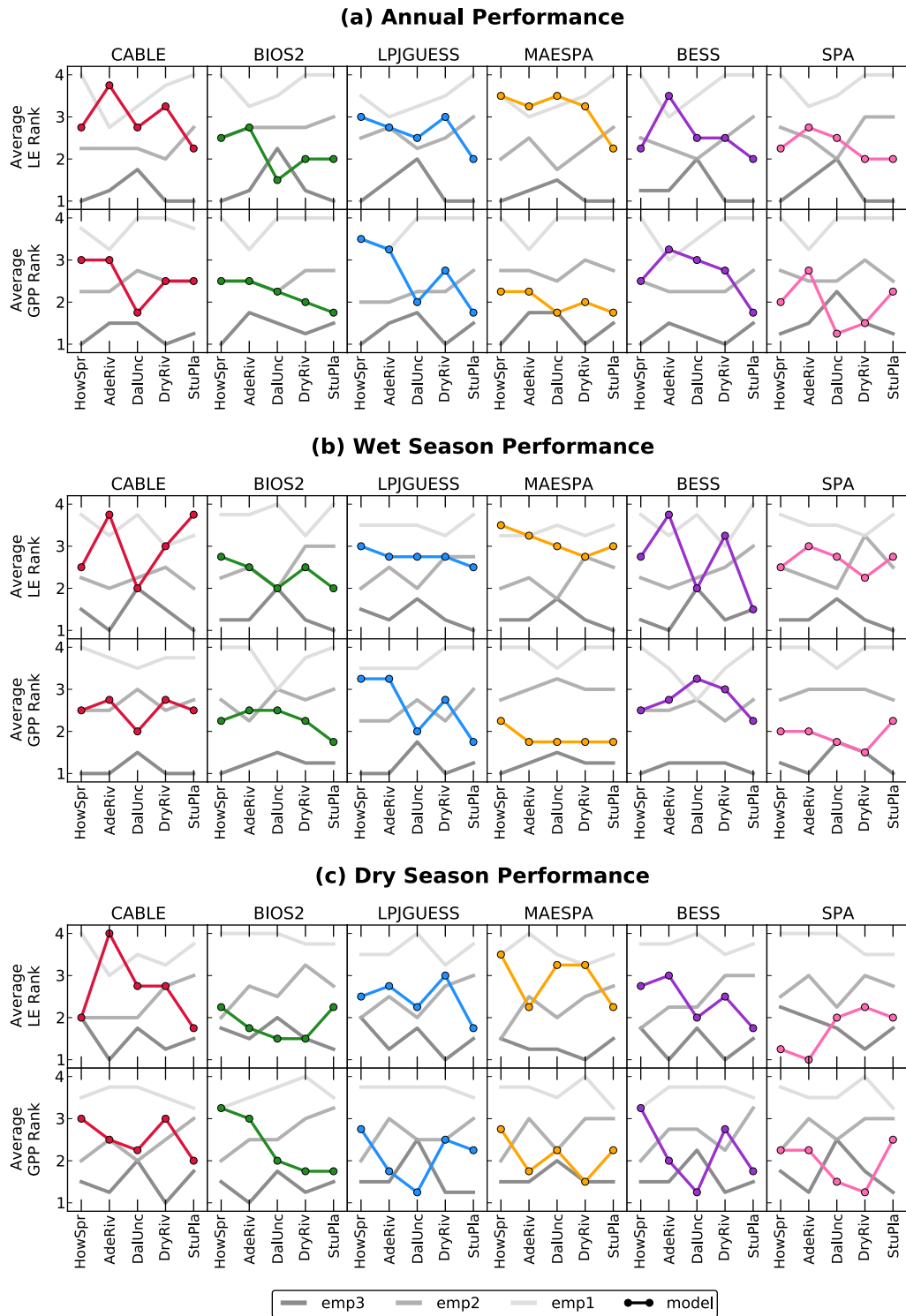


Figure 5: Average rank plot showing the performance of the ecosystem models for all sites across the North Australian Tropical Transect (NATT) ordered in terms of annual rainfall as follows: Howard Springs (HowSpr), Adelaide River (AdrRiv), Daly Uncleared (DalUnc), Dry River (DryRiv), and Sturt Plains (StuPla). Models are individually ranked against the benchmarks in order of 1 to 4 (1 model + 3 benchmarks) and express the amount of metrics the models are meeting in Table B2. The rankings are determined individually for latent energy (LE) and gross primary productivity (GPP). The coloured lines represent each of the 6 models in the study, while the grey lines represent the empirical benchmarks. The average ranking for each model was determined for (a-b) a complete year, (c-d) the wet season and (e-f) the dry season.

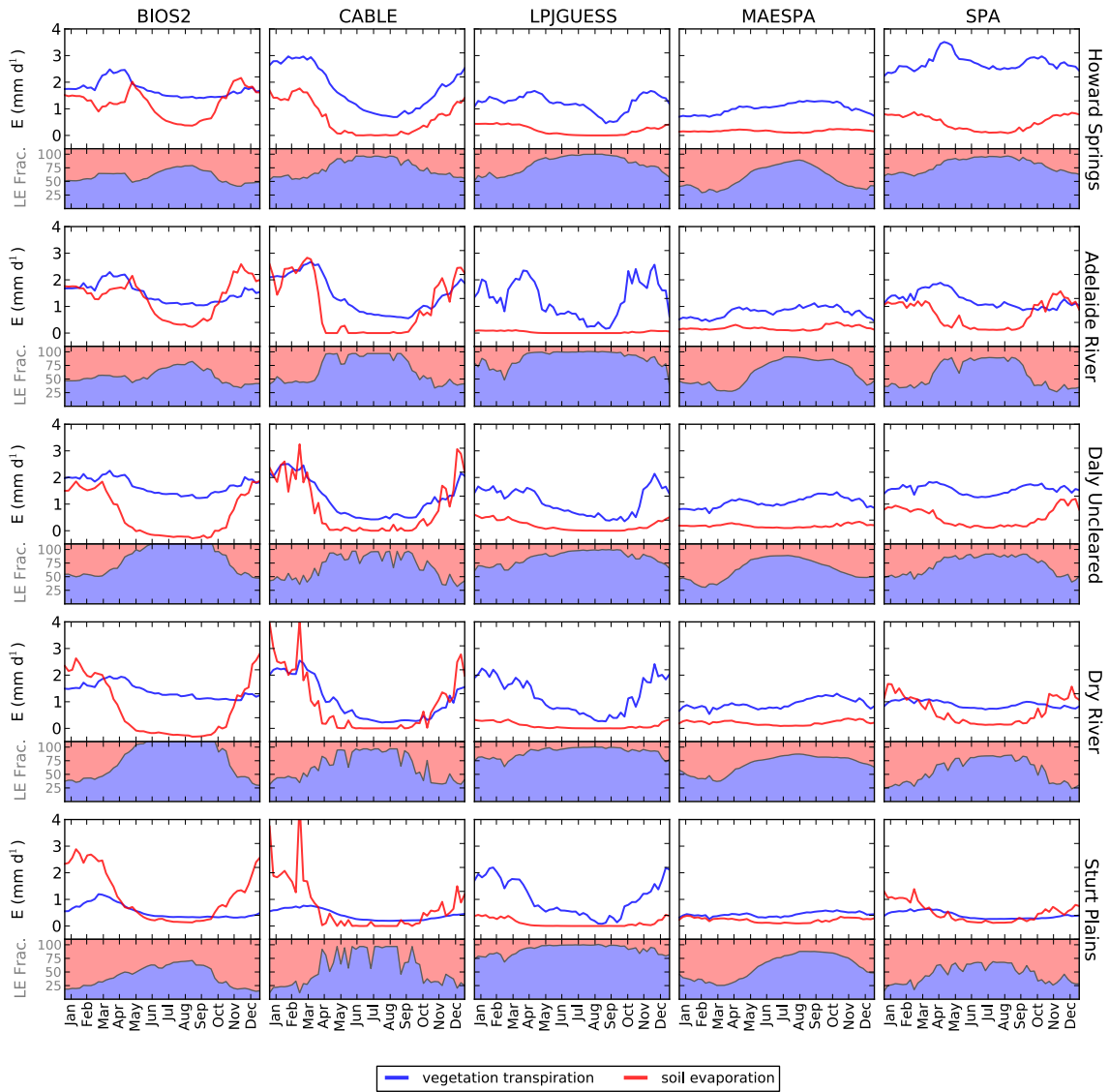


Figure 6: Average year outputs of vegetation transpiration (grass + trees) and soil evaporation, as well as their percentage contributions to total latent energy (LE) for each of the 6 terrestrial biosphere models at each of the 5 savanna sites.

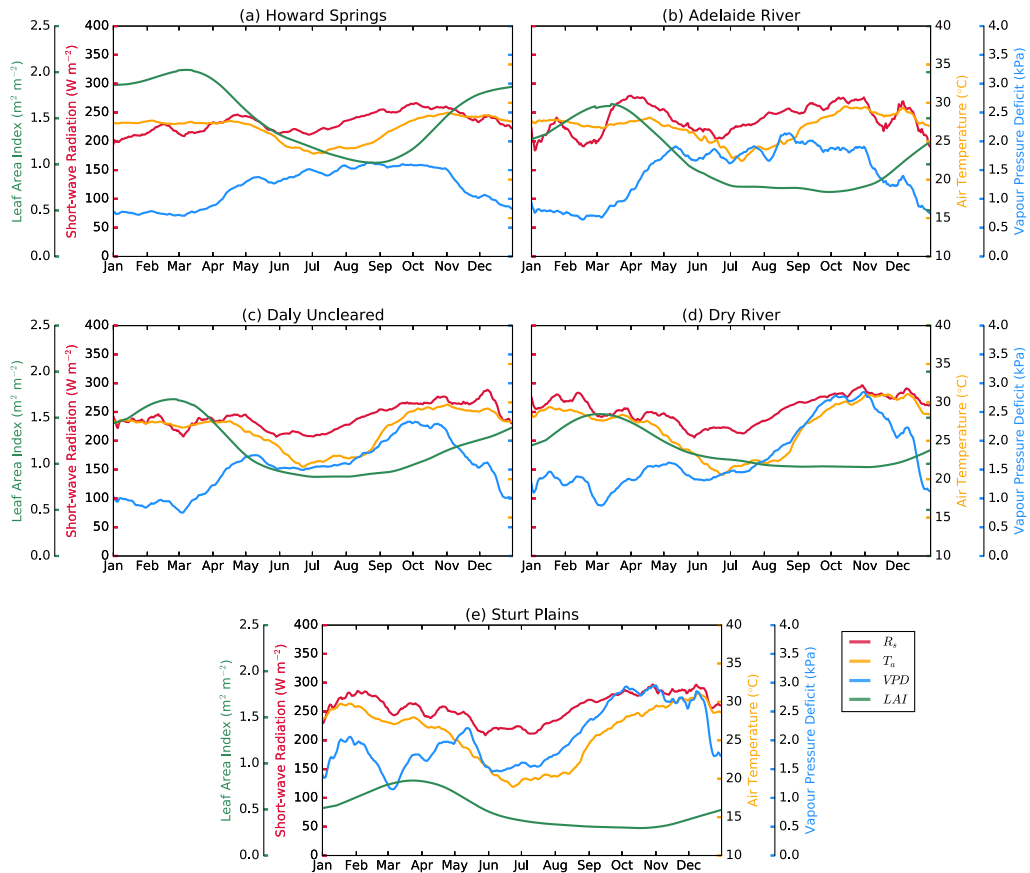


Figure S1: A smoothed (7-day moving average) representation of the environmental drivers used to construct the empirical benchmarks at each of the 5 savanna sites, and are shown from wettest to driest; (a) Howard Springs, (b) Adelaide River, (c) Daly River, (d) Dry River and (e) Sturt Plains. The time-series represents the seasonality over an average year for mean daily solar radiation (R_s), mean daily air temperature (T_a), mean daily vapour pressure deficit (VPD) and leaf area index (LAI).