1      **On the forecast skills of a convection permitting ensemble**

2

3      THERESA SCHELLANDER-GORGAS[1], YONG WANG[1*], FLORIAN MEIER[1],

4      FLORIAN WEIDLE[1], CHRISTOPH WITTMANN[1], ALEXANDER KANN[1]

5

6      [1]Department of forecasting models, Central Institute for Meteorology and

7      Geodynamics, Vienna, Austria

8

9

10

11      (Submitted to Geoscientific Model Development & Discussions)

12

13

14

15

16

17      [*]*Corresponding author address:* Yong Wang, Department of forecasting models,

18      Zentralanstalt für Meteorologie und Geodynamik, Hohe Warte 38, A-1190 Vienna,

19      Austria.

20      E-Mail: yong.wang@zamg.ac.at

21

22

23

24

25                                    ABSTRACT

26

27    The 2.5 km convection-permitting (CP) ensemble AROME-EPS (**A**pplications of

28    **R**esearch to **O**perations at **Me**soscale – **E**nsemble **P**rediction **S**ystem) is evaluated

29    by comparison with the regional 11 km ensemble ALADIN-LAEF (**A**ire **L**imitée

30    **A**daption dynamique **D**éveloppement Inter**N**ational - **L**imited **A**rea **E**nsemble

31    **F**orecasting) to show whether a benefit is provided by a CP EPS. The evaluation

32    focuses on the abilities of the ensembles to quantitatively predict precipitation during

33    a 3-month convective summer period over areas consisting of mountains and

34    lowlands. The statistical verification uses surface observations and 1 km x 1 km

35    precipitation analyses, and the verification scores involve state-of-the-art statistical

36    measures for deterministic and probabilistic forecasts as well as novel spatial

37    verification methods. The results show that the convection-permitting ensemble with

38    higher resolution AROME-EPS outperforms its mesoscale counterpart ALADIN-LAEF

39    for precipitation forecasts. The positive impact is larger for the mountainous areas

40    than for the lowlands. In particular, the diurnal precipitation cycle is improved in

41    AROME-EPS, which leads to a significant improvement of scores at the concerned

42    times of day (up to approximately one third of the scored verification measure).

43    Moreover, there are advantages for higher precipitation thresholds at small spatial

44    scales, which is due to the improved simulation of the spatial structure of

45    precipitation.

46

47

## 1. Introduction

The prediction of deep convection in mountainous terrain is known to be one of the greatest challenges in atmospheric modeling. The initiation and development of deep convection is dependent on small-scale orographic structures and related processes, which cannot be easily described by atmospheric models (Wulfmeyer et al. 2011, Barthlott et al. 2011, Weckwerth et al. 2014). Nevertheless, the estimation of the location, duration and intensity of precipitation events is important as alpine areas are more exposed to natural hazards connected with heavy precipitation (landslides and flooding) than flat land (e.g. Rotach et al. 2009, Haiden et al. 2014).

Models with deep convection-parameterization perform poorly in simulating heavy and highly localized precipitation, especially those with a grid-spacing larger than 10 km (Weusthoff et al. 2010). One source of errors is that the applied convection schemes act independently in individual model grid columns. As a consequence, convectively generated cold-pools that drive convective system propagation cannot be properly simulated, resulting in simulated system movement that is too slow. In weak synoptic forcing, for example, organized MCSs are particularly challenging for convection-parameterizing models (Clark et al. 2007; Liu et al. 2006). Another drawback is that the inadequate descriptions of buoyancy and updrafts in a convection-parameterizing model often cause convection to initiate too early. This premature initiation of convection often results in timing and location errors as well as

68  difficulties to simulate the diurnal cycle of rainfall (Clark et al. 2007). Detailed

69  discussion on the convection initiation in a convection-parameterizing model can be

70  found in Davis et al. (2003) and Bukovsky et al. (2006).

71  A solution for this kind of forecasting problem is offered by a new generation of

72  numerical weather prediction (NWP) models, which have been developed during the

73  last decade. Convection-permitting models with horizontal grid-spacings of

74  approximately 2 km – 3 km offer new possibilities for estimating local impacts. The

75  term *convection permitting* as used in this article (*CP* hereafter) means that a deep

76  convection parameterization is not used in the model. It is assumed that the

77  horizontal resolution around 2-3 km is sufficient to depict the bulk properties of

78  precipitating convective cells, but not to truly resolve the processes within

79  precipitating convective cells such as turbulence and entrainment (Bryan et al. 2003).

80  This is in accordance with Weisman et al. (1997) who suggested setting the upper

81  limit for the range of *convection allowing* resolutions at 4 km.

82  Despite the higher resolution and explicit simulation of deep convection, the exact

83  prediction of location, intensity and spatio-temporal extent of deep convection is still

84  difficult. Recently, probabilistic approaches using convection-permitting ensembles

85  have proven valuable, since they provide direct information on forecast uncertainty,

86  which is often quite large for deep convection. An ensemble usually consists of a

87  number of model runs, which differ in their initial and boundary conditions and/or

88  model configurations. In order to produce a reliable probabilistic forecast, the

89  individual ensemble member forecasts should be equally likely to occur and cover

90  the range of future states. Following Clark et al. (2011), the ideal number of

91  ensemble members is dependent on the point of *diminishing returns,* i.e. the

92  ensemble size where no new information can be expected by additional members.

93  In the recent years several CP EPSs have been developed and some experiences

94  with them have already been made. To name but a few, there are the COSMO-DE-

95  EPS (**Co**nsortium for **S**mall-scale **Mo**deling – EPS, Gebhardt et al. 2011; Peralta et

96  al. 2012; Bouallègue et al. 2013; Kühnlein et al. 2014) at the Deutscher Wetterdienst

97  (DWD), the CP version of UK Met Office's MOGREPS (**M**et **O**ffice **G**lobal and

98  **R**egional **E**nsemble **P**rediction **S**ystem, Bowler et al. 2008; Caron 2013; Hanley et al.

99  2013; Tennant 2015), a Storm Scale Ensemble Forecast (SSEF) run by the Center of

100 Analysis and Prediction of Storms (CAPS) at the University of Oklahoma (Xue et al.

101 2007, 2009; Clark et al. 2011; Schumacher et al. 2013 and Schumacher and Clark

102 2014), WRF based CP ensemble at NCAR (Schwartz et al. 2014) and AROME-EPS

103 (e.g. Vié et al. 2012; Bouttier et al. 2012) developed at Météo-France. A common

104 feature of all of these EPSs is that their horizontal mesh size is equal to or less than

105 4 km, but mostly between 2 km and 3 km.

106 The EPSs mentioned above differ regarding their number of ensemble members and

107 their perturbation strategies and post-processing. Some of them apply an ensemble

108 data assimilation (EDA) approach for perturbing the initial conditions (ICs) (Vié et al.

109 2012; Caron 2013; Schumacher and Clark 2014; Schwartz et al. 2014). The applied

110 model perturbation methods range from a multi-parameter approach (Gebhardt et al.

111 2011) to a stochastic physics scheme (Bouttier et al. 2012; Romine et al. 2014) and

112 to using different dynamical cores (Schumacher et al. 2013). In order to increase

113 ensemble size and to improve the representation of the ensemble distribution some

114   systems also apply the neighborhood method and/or lagged ensemble concepts

115   (Bouallègue et al. 2013). While the neighborhood method is based on ensemble

116   probabilities derived from grid points of a defined environment (Theis et al. 2005,

117   Schwartz et al. 2010), the lagged ensemble approach uses forecasts of successive

118   ensemble runs (Bouallègue et al. 2013).

119   A number of evaluative studies concerned with these CP-EPSs have been

120   conducted. They mainly focus on the investigation of the impact of CP ensemble

121   configurations, for example, the generation of IC perturbation, representation of the

122   model error, uncertainties from the lateral boundary conditions (LBCs), ensemble

123   size, and spatial scale (Kong et al. 2006; Clark et al. 2009; Clark et al. 2011; Vié et

124   al. 2012; Bouttier et al. 2012; Bouallègue et al. 2013; Kühnlein et al. 2014; Schwartz

125   et al. 2014; Schumacher and Clark 2014; Romine et al. 2014; Tennant 2015). There

126   are few comprehensive studies on the evaluation of CP EPS, in particular, in

127   comparison with the mesoscale regional EPS. Clark et al. (2009) compared a 5-

128   member 4 km grid-spacing convection allowing ensemble with a 15-member 20 km

129   grid-spacing regional ensemble. Their case studies reveal that the convection

130   allowing ensemble generally provides more accurate precipitation forecasts than the

131   coarser resolution regional EPS. These results are consistent with those found by

132   Taraphdar et al. (2014) who showed the superior forecast quality of deterministic

133   high-resolution forecasts of tropical cyclone tracks and the accompanying rainfall

134   intensities.

135   In this paper, we will evaluate the performance of a 16-member 2.5 km grid-spacing

136   convection permitting EPS by comparing it with its driving 16--member and 11 km

137  grid-spacing mesoscale regional ensemble. Focus will be on the capabilities of the

138  CP ensemble to quantitatively predict precipitation during a convective summer

139  period over an area consisting of mountains and lowlands. Of interest here is the

140  Alpine region, since the impacts of the mountainous terrain, such as windward/lee

141  effects, the differential heating of valley and mountain slopes can cause large

142  inaccuracies in forecasting convective precipitation and pose a challenge for

143  numerical models and their physical parameterizations (Richard et al. 2007;

144  Wulfmeyer et al. 2008, Bauer et al. 2011, Wulfmeyer et al. 2011). Therefore, an

145  evaluation study is designed and conducted for a typical convective season (3

146  months, May – August 2011), i.e. a period, which is long enough to make at least

147  basic statements about the significance of results. Naturally, this period length is not

148  sufficient to enable statistically reliable statements on real hazardous events, such as

149  landslides and flashfloods. However, the investigations can be regarded as a first

150  step towards this aim. The CP ensemble, which is evaluated in this paper, is a

151  version of AROME-EPS, developed at the Central Institute for Meteorology and

152  Geodynamics in Austria (ZAMG). It is compared with its coarser driving regional EPS

153  ALADIN-LAEF (Wang et al. 2011). The following questions are raised:

154  • Can a convection permitting EPS provide an advantage over its coarser,
155    driving regional EPS in complex terrain?

156  • Is there any difference of the performance for the compared EPSs between
157    lowlands and mountainous areas?

158      • How well can CP EPS and lower resolution regional EPS simulate the diurnal

159           cycle of precipitation? Is the onset and development of convective

160           precipitation realistic?

161      • Does a significant difference in performance for different weather regimes

162           (i.e. days with weak and strong synoptic forcing) exist?

163 A verification study is designed and conducted to answer these questions and to

164 establish whether AROME-EPS can outperform ALADIN-LAEF, a regional

165 mesoscale ensemble with deep convection parameterization on a coarser grid. Wang

166 et al. (2012) demonstrated the added value of ALADIN-LAEF as a regional

167 mesoscale EPS to the global ECMWF-EPS (European Centre for Medium-Range

168 Weather Forecasts). Hence, the present study extends this research by addressing

169 the step between regional mesoscale and CP ensembles.

170 For the present paper, AROME-EPS is coupled to the 16 perturbed ALADIN-LAEF

171 members. This is done to take advantage of the simulation of uncertainties used in

172 ALADIN-LAEF. This uncertainty information is subsequently transferred to finer

173 scales via the dynamical downscaling of the ALADIN-LAEF forecasts by AROME.

174 This means that, both IC perturbations and LBC perturbations are provided from the

175 driving model and are, thus, consistent. No further IC perturbations and model

176 perturbations are applied. Generally, the set-up is kept as simple as possible to point

177 out the pure effects of the downscaling: AROME-EPS is directly coupled to a daily

178 ALADIN-LAEF run initiated at 00 UTC. There is no time lag between the ALADIN-

179 LAEF and the AROME-EPS simulations and the forecasts are evaluated for the first

180 30h of the model runs, hence for a whole day and the subsequent night each.

181 The benefits of AROME-EPS compared to ALADIN-LAEF are revealed in the

182 framework of a comparative verification study. Although the focus of the verification

183 study is on the onset and development of precipitation, the performance of other

184 surface weather parameters are considered. The verification methods are selected in

185 such a way that the overall performance, in a deterministic and probabilistic manner,

186 and the abilities of the ensembles to reproduce spatial structures, can be

187 investigated. Hence, ensemble-related scores are combined with spatial verification

188 methods.

189 More detailed characteristics of the compared models are described in Section 2

190 along with the verification data. The methods chosen for the evaluation of the two

191 ensembles are described in Section 3. Section 4 comprises the verification results

192 and Section 5 the summary and concluding remarks.

193 **2. Ensemble systems and data**

194 *a. The regional ensemble ALADIN-LAEF*

195 ALADIN-LAEF is the operational regional ensemble system of ZAMG and runs at

196 ECMWF (Wang et al. 2010, 2011). It is based on the hydrostatic spectral limited area

197 model ALADIN (Wang et al. 2009). ALADIN-LAEF has 16 members and is coupled to

198 ECMWF-EPS (Weidle et al. 2013) with a horizontal grid-spacing of 11 km. In

199 operational mode it and runs two times per day at 0000 and 1200 UTC and provides

200 probabilistic forecasts on a forecast range up to 3 days ahead, i.e. 72 h. In this study,

201 however, evaluation is confined to the run at 00 UTC and a forecast range of 30 h

202    ahead only. This is done in order to investigate the onset and development of
203    convection in its diurnal cycle.

204    ~~with a horizontal grid-spacing of 11 km.~~ The 16 members of ALADIN-LAEF are not
205    sufficient to represent the atmospheric state probability density function (PDF).
206    However, Schwartz et al. (2014) have shown that similar verification scores can be
207    obtained from a 50-member ensemble and subsets of 20-30 members. Hence, we
208    can expect, at least, reasonable results from verification based on a 16-member
209    ensemble.

210    ~~The goal of ALADIN-LAEF is to provide probabilistic forecasts on a forecast range up~~
211    ~~to 3 days ahead, i.e. 72 h, although only 30 h are used in this study for the~~
212    ~~comparison with AROME-EPS.~~ The ALADIN-LAEF domain (Figure 1) covers the
213    whole European continent, Iceland, the whole Mediterranean Sea, Black Sea,
214    Caspian Sea and adjacent countries. The eastern margins reach the Ural Mountains
215    and parts of Siberia. To deal with the atmospheric initial condition perturbation
216    ALADIN-LAEF applies a breeding-blending method for generating the IC
217    perturbations for the upper levels. It uses large-scale perturbations from the driving
218    global-ECMWF-EPS combined with small-scale perturbations from the ALADIN-
219    breeding vectors (Toth and Kalnay 1993). The blending method (Wang et al. 2014)
220    ensures that inconsistencies between small and large-scale perturbations are
221    avoided. Therefore a digital filter is applied on the low spectral truncations of both the
222    breeding-vectors and the fields from the global model. Afterwards the filtered
223    breeding vectors on the full spectral resolution are subtracted from the original ones

224  and added by the filtered global fields resulting in initial perturbations that are

225  consistent with the regional EPS itself as well as with the driving global EPS.

226  To consider uncertainties arising from the initial surface conditions in ALADIN-LAEF,

227  a surface data assimilation scheme based on optimum interpolation (CANARI - Code

228  for the Analysis Necessary for Arpege for its Rejects and its Initialization, Taillefer

229  2002) is implemented using randomly perturbed observations. To account for

230  uncertainties in the model itself, a multi-physics approach is implemented in ALADIN-

231  LAEF. The perturbed members use different model configurations with several

232  combinations and tunings of schemes and parameterizations available in the ALADIN

233  physics package. The main emphasis is put on the variation and tunings of the

234  following schemes and parameterizations: The diagnostic convection scheme as

235  described in Bougeault (1985); the prognostic deep convection scheme 3MT

236  (modular multi-scale **M**icrophysics and **T**ransport scheme, Gerard et al. 2009), and

237  the connected microphysics scheme described in Geleyn et al. 2008 and Gerard et

238  al. (2009); the radiation scheme based on Ritter and Geleyn (1992) or alternatively

239  the scheme described in Mlawer (1997) and Morcrette (1991); the pseudo prognostic

240  TKE (**T**urbulent **K**inetic **E**nergy) scheme described in Vana et al. (2008). Further

241  details can be found in (Wang et al. 2010).

242  *b. The convection permitting ensemble AROME-EPS*

243  The model core of AROME-EPS is the non-hydrostatic spectral limited area model

244  AROME (Seity et al. 2011), which is especially designed to run at very high

245  resolutions with a grid-spacing of 2.5 km or lower. Deep convection is treated

246  explicitly, while shallow convection is parameterized with a mass flux approach

247 (Pergaud et al. 2009). The single moment bulk microphysics scheme ICE3 for mixed-

248 phase cloud parameterization (Pinty and Jabouille 1998) can handle mixing ratios of

249 five prognostic hydrometeor classes: cloud water, cloud ice, rain, snow and graupel

250 and also simulates complex interactions between them. AROME by default uses a

251 three-layer soil model SURFEX (Surface Externalisé) with the effects of sea and

252 urban areas parameterized using a tile approach (Masson et al. 2000).

253 At ZAMG a deterministic version of AROME with 2.5 km grid-spacing has been

254 operational since January 2014 running every 3 hours up to a lead-time of 48 hours.

255 The domain for the model integration encompasses the Alpine region (Figure 1).

256 Table 1 summarizes the most important model characteristics of ALADIN-LAEF and

257 AROME-EPS.

258 To run AROME-EPS, the same version of AROME with the same resolution is

259 initialized by a dynamical downscaling of ALADIN-LAEF and coupled to the 16

260 members of ALADIN-LAEF. The ensemble runs with a forecast range of 30 h are

261 initiated at 00 UTC each day, i.e. at the same time as ALADIN-LAEF. There is no A

262 time lag  is not considered, as the pure impact of enhanced resolution and the

263 convection-permitting configuration shall be investigated. Apart from the

264 perturbations of initial conditions and lateral boundary conditions, no further

265 perturbations (such as e.g. multi-physics parameterizations as in ALADIN-LAEF) are

266 induced in the model integration. This comparatively simple configuration is used for

267 several reasons: First, AROME-EPS has been set up quite recently at ZAMG and is

268 still at an early stage of development. Secondly, the development of physics

269 perturbations in AROME-EPS will rather go towards a stochastic physics scheme or

270 a combined stochastic/multi-physics scheme than towards pure multi-physics as

271 currently used in ALADIN-LAEF. And thirdly, the aim of this study is to test the

272 possible advantage of a CP EPS compared to the operational system of ALADIN-

273 LAEF.

274 *c. Verification data*

275 Station observations are used for the evaluation of ALADIN-LAEF and AROME-EPS

276 surface weather variables. Figure 2 shows the 517 surface stations in the AROME

277 domain, providing observations at 6-hourly intervals for 2 m temperature, 2 m

278 humidity, 10 m wind speed and mean sea level pressure. The upper level verification

279 is achieved using ECMWF analyses reference data at four pressure levels: 925 hPa,

280 850 hPa, 700 hPa, and 500 hPa, which are adapted to the model resolutions of both

281 AROME-EPS and ALADIN-LAEF.

282 The evaluation of precipitation forecasts is performed using the very high-resolution

283 precipitation analyses of the ZAMG nowcasting system INCA (Integrated Nowcasting

284 through Comprehensive Analyses; Haiden et al. 2011). This is necessary as the

285 average station distance of precipitation observations is too large to resolve the fine

286 spatial structures of precipitation events. The advantage of the INCA analyses is that

287 they use additional observations and are provided on a regular grid. Based on this

288 gridded data, it is possible to apply enhanced verification methods on precipitation

289 fields, which cannot be computed on a point-to-point basis.

290 The INCA system, developed at ZAMG, operates on a horizontal resolution of 1 km x

291 1 km. INCA blends data from automatic weather stations, remote sensing data

292  (radar, satellite), forecast fields of numerical weather prediction (NWP) models, and

293  high-resolution topographic data (Haiden et al. 2011). It provides hourly 3-D fields of

294  temperature, humidity, wind, and 2-D fields of cloudiness, precipitation rate and

295  precipitation type with an update frequency of 15 minutes to 1 hour. The precipitation

296  analyses are provided for different accumulation periods. In the present study, the

297  one-hour accumulated INCA precipitation analyses are used as a reference for the

298  spatial verification of EPS forecasts. For these analyses, precipitation measurements

299  from surface stations and radar data are accumulated to one-hour sums and

300  algorithmically merged. Prior to the analysis procedure, the data are quality

301  controlled and climatologically scaled (Haiden et al. 2011). In this way the higher

302  quantitative accuracy of the station data and the better spatial coverage of the radar

303  data are utilized. The resulting analysis reproduces the observed values at the

304  station locations while preserving the spatial structure provided by the radar data.

305  The analysis error, which is computed from classical cross-validation, varies from

306  case to case and depends on precipitation type, e.g. large-scale or convective, and

307  on the accumulation period. The magnitude of analysis errors of grid point values can

308  be quite large, but areal mean values are significantly more reliable (Haiden et al.

309  2011)

310  Amending the rain gauge - radar combination, the scheme includes elevation effects

311  on precipitation using an intensity-dependent parameterization (Haiden and Pistotnik

312  2009). A NWP model first guess is not required in the precipitation analysis, thus

313  such analyses are ideally suited as an independent reference to validate NWP

314  models.

315 Forecast verifications are performed at the observation locations for surface variables

316 as 2 m temperature and humidity, 10 m wind speed and mean sea level pressure,

317 and on the INCA grid for precipitation. The model forecasts are interpolated bi-

318 linearly to the station locations and INCA analysis grid points, respectively. Further, a

319 height correction scheme is applied on 2 m temperature values based on

320 atmospheric standard conditions. In doing so, the same number of

321 forecast/observations pairs is available for the verification of each of the EPS models.

322 This supports the comparability of the verification results.

323 **3. Verification strategy**

324 AROME-EPS and ALADIN-LAEF are evaluated over a 3-month summer period from

325 15 May, 2011 – 15 August, 2011, which represents a typical convective summer

326 season in Central Europe.

327 Precipitation is one of the parameters for which the biggest improvement is expected

328 from the convection-permitting models. Therefore, the evaluation of the ensembles

329 focuses on the representation of the spatio-temporal structure of precipitation events

330 in the forecasts. Nevertheless, the preconditions for the development and onset of

331 precipitation are also considered. For this reason other forecast parameters, such as

332 temperature, humidity, wind speed, air pressure and geopotential height are also

333 verified.

334 Precipitation forecasts are evaluated in both deterministic and probabilistic ways. The

335 deterministic approach is directed towards predicting the correct precipitation

336 amounts and the spatial distribution of the data. Probabilistic evaluation tests the

337 capability of the ensembles to predict a pre-defined event with the probability, which
338 corresponds to its relative frequency, i.e. to produce a reliable PDF for the
339 occurrence of the event. The events can be defined as, e.g., precipitation amounts
340 exceeding a certain threshold. In this study, thresholds of 0.1 mm (threshold for the
341 prediction of *rain* or *no rain*), 0.5 mm, 1 mm, 2 mm and 5 mm are chosen for 3-hourly
342 accumulated precipitation amounts. These thresholds appear low, especially when
343 taking into account convective precipitation events. However, the thresholds are
344 selected according to the frequency of occurrence of the precipitation values in the
345 individual grid cells of the 1 km x 1 km verification grid. They ensure that a sufficient
346 number of observed events are available for evaluation over the 3-month test period.
347 The two ways of deterministic and probabilistic evaluation reflect the main options for
348 the efficient use of ensemble forecasts: First, as a conservative prediction of
349 ensemble mean or median or, second, as a tool to estimate the uncertainty of the
350 forecast and the probability of extreme values via the ensemble spread and PDF
351 (e.g. Zhu et al. 2002).

352 A number of t~~T~~raditional point-to-point verification scores (see e.g. Wilks 2006) ~~in~~
353 ~~Table 2~~ are computed for all evaluated parameters. In addition, significance tests for
354 these scores are performed. Confidence intervals of the verification scores are
355 estimated by a bootstrapping algorithm (Davison and Hinkley 1997; Joliffe 2007;
356 Ferro 2007) and confidence intervals of 90%. The bootstrapping method uses 5000
357 random samples with a block length of eight.

358 In order to present the results concisely, ~~only~~ three scores have been selected ~~from~~
359 ~~Table 2~~ to describe the differences in forecast performance between AROME-EPS

360    and ALADIN-LAEF: Bias (Eq. 1), Brier Score (BS, Brier 1950, Eq. 2) and Continuous

361    Ranked Probability Score (CRPS, Hersbach 2000; Gneiting and Raftery 2007; Eq. 3).

362    The Bias simply measures the mean deviation between the analyzed values (*a)* and

363    the forecast values, in our case the ensemble means $\left(\overline{f}\right)$, at *n* grid points *i*. Both,

364    positive as well as negative signs are possible. A perfect forecast has a bias of zero.

365    (1)        $$Bias = \frac{1}{n}\sum_{i=1}^{n}\left(\overline{f}_i - a_i\right)$$

366    Like the Bias also BS is a measure for the accuracy of the forecasts, however, in

367    probability space. It is the mean squared difference between the forecast probability

368    $p$  ($p \in \left[0;1\right]$, e.g. derived from the distribution of ensemble members) for a pre-

369    defined event (e.g. the exceeding of a threshold) and the analyzed truth $x$  ($x \in \left\{0,1\right\}$

370    ). The minimal value of zero is achieved for a perfect forecast, and the maximum

371    value is one for the worst possible forecast.

372    (2)        $$BS = \frac{1}{n}\sum_{i=1}^{n}(p_i - x_i)^2$$

373    CRPS is related to BS insofar, as it can be expressed as the integral of BS for all

374    possible thresholds of the meteorological parameter $\xi$  (Hersbach 2000). The value

375    for an ideal forecast of CRPS is zero as for BS.

376    (3)        $$CRPS_i = \int_{-\infty}^{\infty}\left[P_i\left(\xi\right) - P_i\left(\xi_a\right)\right]^2 d\xi$$

377 The continuous ranked probability score compares the cumulative distributions $P_i(\xi)$

378 (Eq. 4) and $P_i(\xi_a)$ (Eq. 5) of the forecast and the analyzed values at each grid point $i$

379 .

380 (4)
$$P_i(\xi) = \int_{-\infty}^{\xi} p_i(y)\,dy$$

381 (5)
$$P_i(\xi_a) = H(\xi - \xi_a)$$

382 $H(\xi)$ is the so-called Heaviside-function (Eq. 6), which only takes the values 0 and 1.

383 (6)
$$H(\xi) = \begin{cases} 0 & for \quad \xi < 0 \\ 1 & for \quad \xi \geq 0 \end{cases}$$

384 In addition to those traditional statistical scores ~~in Table 2~~, precipitation forecasts are

385 verified by spatial verification methods, which not only consider the exact match of

386 forecast and verification values at individual points, but take into account the

387 matching of forecasts and observations in terms of objects or spatial scales (Casati

388 et al. 2008, Ahijevych et al. 2009, Gilleland et al. 2010). This is necessary as

389 precipitation fields exhibit high spatial variability and discontinuity. Small deviations in

390 space and time between forecast and verification data can lead to large errors in

391 traditional point to point verification scores, which is also known as the *double*

392 *penalty* problem (Nurmi 2003).

393 *a. Spatial verification methods*

394    The selected spatial verification methods are the so-called SAL method (Structure-

395    Amplitude-Location method, Wernli et al. 2008) and the Fractions Skill Score

396    (Roberts and Lean 2008).

397    SAL determines the forecast performance in terms of structure (S), amplitude (A) and

398    location (L). The method is object based. Precipitation objects in forecast and

399    verification fields are contiguous areas of grid-points exceeding a certain precipitation

400    threshold.

401    (7)            $$A = \frac{\bar{R}_f - \bar{R}_a}{0.5\left[\bar{R}_f + \bar{R}_a\right]}$$

402    The amplitude score (Eq. 7) defines whether the integrated precipitation amount $\bar{R}$ of

403    the field is underestimated (A < 0) or overestimated (A > 0). Subscripts, $f$ and $a,$

404    denote forecast and analyzed fields, respectively.

405    The location score measures the agreement of the centers of mass in the analyzed

406    and predicted precipitation fields together with the averaged distance between the

407    center of mass and the individual objects. It is actually the sum of two components L=

408    L1+ L2 where both values are in the range [*0, 1*]. The first part L1

409    (8)            $$L1 = \frac{\left|x(R_f) - x(R_a)\right|}{d_{max}}$$ ,

410    is a measure of the distance between the mass centers $x$ of the analyzed ($R_a$) and

411    the predicted precipitation fields ($R_f$). $d_{max}$ is the longest possible distance in the

412    domain.

413  As an identical mass center position does not necessarily mean that the forecast is

414  perfect, the second component L2 (Eq. 9) is introduced:

415  (9)
$$L2 = 2\frac{\left|r(R_f) - r(R_a)\right|}{d_{max}}.$$

416  L2 takes into account the distance $r$ (Eq. 10) between the mass center of each

417  individual object $R_n$ and the overall mass center and compared between the

418  observed and simulated precipitation field:

419  (10)
$$r = \frac{\sum_{n=1}^{M} R_n \left|x - x_n\right|}{\sum_{n=1}^{M} R_n}.$$

420  The L component has a range [0, 2] with $L=0$ indicating a perfect forecast.

421  The structure score $S$

422  (11)
$$S = \frac{V(R_f) - V(R_a)}{0.5\left[V(R_f) + V(R_a)\right]}$$

423  compares the weighted sums of the precipitation volumes $V(R)$

424  (12)
$$V(R) = \frac{\sum_{n=1}^{M} R_n V_n}{\sum_{n=1}^{M} R_n}$$

425  of the precipitation objects, where the $V_n = R_n / R_{max}$ describe precipitation sums

426  scaled by their maxima. If S < 0, forecast objects are too small and too peaked. In

427  contrast, S > 0 indicates that the objects are too large and too flat.

428   The fractions skill score (FSS)

429   (13)
$$FSS(n) = 1 - \frac{MSE(n)}{MSE(n)_{ref}}$$

430   evaluates the forecasts on different spatial scales. The scales are defined via

431   neighborhoods, i.e. square boxes of length $n$ grid spaces surrounding a selected grid

432   point. The score compares the fractions of rain coverage of forecast and analysis in

433   the neighborhoods. Depending on the precipitation event, small disparities of the

434   coverage may lead to large forecast errors on fine scales, but to a better rating on a

435   coarser scale. The aim of FSS is to identify scales for which the evaluated model can

436   provide useful forecasts.

437   FSS is computed by assigning the grid points binary values 0 and 1 in each of the

438   neighborhoods with subscripts $(i,j)$, according to a selected precipitation threshold.

439   From these binary fields, the fraction of the points with value 1 are computed for

440   analyses and forecasts as $A_{(n)i,j}$ and $F_{(n)i,j}$ , respectively.

441   At each such defined scale $n$, the mean squared error ($MSE$):

442   (14)
$$MSE_{(n)} = \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} \left[ A_{(n)i,j} - F_{(n)i,j} \right]^2$$

443   is computed for the whole field of fractions and related to a reference ($MSE_{ref}$)

444

445     (15)

$$MSE_{(n)ref} = \frac{1}{N_x N_y} \left[ \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} A^2_{(n)i,j} + \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} F^2_{(n)i,j} \right].$$

446     $MSE_{ref}$ is the largest possible MSE, which can be obtained from the underlying field.

447     The skill score summarizes the performance in the whole field and ranges from 0

448     (complete mismatch) to 1 (perfect match).

449     *b. Subdomains for precipitation verification*

450     Verification is done for the whole domain *Austria*. To account for the different

451     topographic characteristics in the verification domain, two sub-domains are chosen

452     (Figure 3). They comprise mountainous area (region *West*) as well as region with flat

453     terrain (region *Northeast*). Due to the location of the Alps in Austria and the prevailing

454     flow directions around the Alps, each of the subdomains has its own climatological

455     properties, which is also visible in the precipitation characteristics.

456     *c. Temporal stratification*

457     In order to investigate the influence of different weather regimes, the 92 days of the

458     test period are classified into three bins according to the synoptic situation*, strong*

459     *synoptic forcing*, *weak synoptic forcing,* and *dry*. Days are classified as *dry* (5 days) if

460     the areal mean of the daily precipitation sum is below 0.05 mm. All other days, i.e. 87

461     days on which rains was reported, are assigned to the bins of *weak (*23 days) or

462     *strong synoptic forcing* (64 days). For the classification, a method described by Done

463     et al. (2006) and successfully applied by Kühnlein et al. (2014) is used which is

464     based on the temporal variability of CAPE (**C**onvective **A**vailable **P**otential **E**nergy) as

465     a measure of atmospheric instability. According to Done et al. (2006), the approach

466 helps to distinguish between days on which convection is predominantly at

467 *equilibrium* or at *non-equilibrium.* This means that the destabilization of the

468 atmosphere by large-scale synoptic forcing is balanced or un-balanced, respectively,

469 by the stabilization through convection. The idea is that this balance or imbalance is

470 related to the timescale in which CAPE is built up by large-scale processes and

471 consumed by convection. On days with *weak synoptic forcing* the consumption of

472 CAPE is related to the diurnal cycle or to local triggering rather than to prevalent

473 large-scale processes. In these cases the convective timescale is long and CAPE is

474 often not fully consumed by convection. In situations where CAPE is realized much

475 faster by large-scale processes, i.e. in situations of *strong synoptic forcing,*

476 convection is in equilibrium. In our study the *convective adjustment time-scale* $t_c$

477 (16) $$t_c = CAPE \frac{d(CAPE)}{dt}$$

478 is calculated hourly from AROME-EPS CAPE forecasts using $\Delta t = 1h$. Following the

479 suggestion of Done et al. (2006) a specific day is assigned to *weak synoptic forcing* if

480 the areal mean of $t_c$ exceeds a threshold of 6 h at least once a day by at least three

481 ensemble members. In order to test the method of Done et al. (2006) we compared

482 the classification with alternative approaches, such as the temporal change of mid-

483 tropospheric vorticity and convection related to patterns in 500 hPa geopotential

484 using archived ECMWF forecast and ERA-Interim re-analyses. The results were

485 comparable to those of the equilibrium method.

486 **4. Results**

23

487 In the following we present the evaluation of AROME-EPS and ALADIN-LAEF over a

488 three-month summer period. The focus is on the performance of near surface

489 parameters, in particular the precipitation forecast, which is of most interest to the

490 users of convection permitting and regional EPSs.

491 *a. Evaluation of forecasts of temperature, wind and humidity*

492 The forecast performance of surface parameters (2 m temperature and humidity, 10

493 m wind speed and mean sea level pressure MSLP) and upper level parameters

494 (temperature, humidity, wind speed and geopotential height) of AROME-EPS and

495 ALADIN-LAEF are verified in this study, which form the background of the evaluation

496 of precipitation.

497 A large number of verification metrics (Table 2) have been calculated for those near

498 surface and upper air parameters. In general there is no clear advantage either for

499 ALADIN-LAEF or for AROME-EPS. Exceptions from this statement are solely

500 constituted by biases in the forecasts, which are particularly found on the surface

501 level. They form the most eminent differences in the performances of the EPSs: If the

502 bias is low, the models provide good performance also for other scores.

503 For the surface level, we also found more results on a high level of significance (i.e.

504 90%). The verification results of the upper levels are less significant than for the

505 surface and performance is more ambivalent. We used a large number of

506 observations for both surface (station observations) and upper levels (ECMWF grid

507 values). Hence, the lower significance of the results for the upper levels can be

508 explained by the model set-up rather than by the verification data. Near surface and

509  on lower levels AROME-EPS can add more information to the model simulation

510  compared to ALADIN-LAEF than on higher levels. This is due to the SURFEX soil

511  scheme and the interaction between a refined representation of orography and the

512  model physics schemes and dynamics. On the higher levels, however, there is less

513  influence of the orography and the simulation resembles more the driving model. For

514  this reasonTherefore, surface results have been selected to highlight the main

515  findings in the following.

516  Figure 4 compares the bias and Continuous Ranked Probability Score (CRPS, see

517  Wilks 2006 for details) for 2 m relative humidity, 2 m temperature and 10 m wind

518  speed. CRPS compares the forecast PDF to the observed values of occurrence and

519  non-occurrence, respectively. CRPS is sensitive to the difference between the

520  forecast probabilities to observed values. The lower the difference, the better the

521  forecast is rated. Hence, the value of CRPS of a perfect forecast is zero. Due to its

522  formulation, signals of CRPS are also reflected by many other scores, in particular

523  those which are sensitive to deviations between the distributions of forecasts and

524  observations. Thus, CRPS is useful for representing the results of this study

525  exemplarily. It also shows the impact of biased forecasts.

526  Biases of 2 m relative humidity in Fig. 4a show noticeable diurnal variations. During

527  the night and early morning, AROME-EPS is too dry, whereas ALADIN-LAEF is too

528  moist during the day (1200 UTC and 1800 UTC). The diurnal variations of the

529  differences between AROME-EPS and ALADIN-LAEF are also reflected in CRPS in

530  Figure 4b. During the night, AROME-EPS and ALADIN-LAEF are at the same level,

531  but for the day hours AROME-EPS shows better results. For 2 m relative humidity,

25

532    most verification results are significant at a level of 90%. This is also true for the

533    differences in forecast performance during the day hours. Results for 2 m

534    temperature in Figures 4c and 4d show an improvement for most of the used scores

535    at a significance level of 90% for AROME-EPS. This result is partially due to a large

536    bias of ALADIN-LAEF temperatures. In contrast, there exist fewer deviations

537    between the ensembles for wind speed (Figures 4e and 4f) and MSLP (not shown).

538    However, these results have only a low level of significance.

539    *b. Evaluation of precipitation forecasts*

540    Precipitation is evaluated by 3-hourly INCA analyses on a regular 1 km x 1 km grid. A

541    first insight of the strengths and weaknesses of the ensembles in forecasting

542    precipitation is offered by a comparison of the daily variability of precipitation

543    intensities. Figure 5 compares the 3-hourly precipitation sums of INCA and both EPS

544    models for different regional domains and for days with strong (left panels) and weak

545    (right panels) synoptic forcing.

546    Errors occur in terms of over- and underestimation of the maximum intensity and in

547    terms of time shifts. The daily maximum of 3 h-precipitation is overestimated by

548    AROME-EPS for regions *West* and *Austria* and both types of synoptic forcing by

549    20%-50%. In ALADIN-LAEF, the maximum in these regions is approximately at the

550    same level as analyzed by INCA. Hence, the too moist conditions of ALADIN-LAEF

551    near the surface in Fig. 4a are not reflected in the precipitation sums. For region

552    *Northeast,* AROME-EPS correctly simulates the maximum amount of precipitation,

553    whereas ALADIN-LAEF is too low.

554    Considering the days with strong synoptic forcing in Figure 5 (left panels), the highest

555    precipitation sums are detected around 1800 UTC. AROME-EPS describes the

556    temporal maximum quite well, whereas the maximum in ALADIN-LAEF occurs too

557    early (-3 h time shift).  In the case of weak synoptic forcing shown in Figure 5 (right

558    panels), the precipitation maxima are observed later than for the other cases in

559    region *West* (e.g. 2100 UTC instead of 1800 UTC). This is not reflected by the EPS

560    models, which both reach the maximum intensity of precipitation at 1500 UTC. Only

561    for region *Northeast* and weak synoptic forcing does the maximum of precipitation

562    occur too late in AROME-EPS. The characteristic that ALADIN-LAEF and AROME-

563    EPS tend to trigger moist and deep convection over complex orography too early is

564    well known (Wittmann et al. 2010). However, according to Figure 5, running a model

565    or an EPS on CP scales is beneficial for predicting the daily maximum of the

566    convective diurnal cycle, at least over mountainous terrain. With respect to the timing

567    of the maxima, AROME-EPS shows a time shift of -3 h, with ALADIN-LAEF -6 h for

568    weak synoptic forcing in regions *Austria* and *West* (panels b) and d), respectively).

569    Because of the limited framework of this study we can only speculate that this

570    behavior might be due to differences caused by the deep convection scheme in

571    ALADIN-LAEF, which is one of the reasons to cause an early onset of precipitation

572    (Bechtold et al. 2013), and respectively, the explicit simulation of deep convection in

573    AROME. Another reason, which we cannot exclude, could be that ALADIN-LAEF and

574    AROME apply different physical parameterizations. The different dynamical cores,

575    hydrostatic and non-hydrostatic, might also contribute to the differences to some

576    extent, but remain statistically less significant in respect of precipitation as shown in

577    an earlier study (Wittmann et al. 2010). Experiences concerning the pure impact of

578 different vertical resolutions on the forecast quality are few. However, it is known that

579 an increase of vertical resolution and, hence, enhanced possibilities to simulate

580 convection-related, micro-physical and boundary-layer processes, does not

581 necessarily result in an improvement of precipitation forecasts. It is rather related to

582 increased overprediction of precipitation amounts (Aligo et al. 2009).

583 A further characteristic evident in Figure 5, is that the precipitation amounts in

584 AROME-EPS develop independently of those in the driving ALADIN-LAEF members,

585 which is indicated by the ensemble spread. In ALADIN-LAEF the ensemble spread is

586 quite large for certain lead times, ranging from a larger overestimation of the

587 observed precipitation amounts to a large underestimation. This contrasts with

588 AROME-EPS, which shows a much smaller range of precipitation amounts. This

589 difference in the spread is very likely due to the large influence of the multi-physics

590 configuration in ALADIN-LAEF, compared with the single physics configuration of

591 AROME-EPS. The scores, which are discussed in the following, Brier score, SAL

592 scores and fractions skill score, demonstrate in which ways the differences in the

593 diurnal precipitation cycle have an influence on forecast quality.

594 *i. Brier score*

595 Figure 6 shows the differences of the Brier Score (BS; Brier 1950), for strong and

596 weak synoptic forcing with different precipitation thresholds. BS measures the

597 accuracy of probability forecasts, which is equivalent to the MSE for deterministic

598 forecasts. The value for perfect forecasts is zero. BS has largest values for the

599 lowest precipitation threshold (0.1 mm, upper panels), and decreases for larger

600 thresholds (2 mm, lower panels).

601  During the morning hours (+6 h, +30 h lead time), BS is low for days with weak

602  synoptic forcing. This is due to the fact, that on these days, generally stable

603  conditions prevail in the morning and precipitation probability is very low. For the

604  lower precipitation threshold, AROME-EPS shows significantly better values than

605  ALADIN-LAEF from 0900 UTC to 1500 UTC. This applies for both, days with weak

606  synoptic forcing and days with strong synoptic forcing.

607  The differences in BS between ALADIN-LAEF and AROME-EPS can, for the most

608  part, be explained by the fact that the precipitation generally starts too early in

609  ALADIN-LAEF forecasts. Additionally, the tendency of ALADIN-LAEF to forecast

610  smoother precipitation fields than AROME-EPS can be assumed as a second source

611  of errors. The smoothness leads to rather medium precipitation probabilities in large

612  areas. BS, however, accounts for sharp forecasts near zero and one (i.e. very low

613  and very high probabilities for rainfall).

614  *ii. SAL scores*

615  The variability of SAL scores with lead-time gives insight in the performance of

616  AROME-EPS and ALADIN-LAEF in terms of the structure, amplitude, and location of

617  the predicted precipitation events. Figures 7 and 8 show the SAL scores for the

618  mountainous region *West* and the lowland region *Northeast,* respectively. The

619  distributions of SAL values are sampled for the individual ensemble members and

620  classified into days with strong (panels a and b) and weak synoptic forcing (panels c

621  and d). These values differ from those based on the ensemble mean and median

622  forecasts as the averaging produces more smoothed precipitation events and, hence,

623  has an influence on the properties described by the SAL method.

624    In both geographic regions and for both types of synoptic forcing, the structure score

625    is lower for AROME-EPS than for ALADIN-LAEF, which is, inter alia, a consequence

626    of the model resolution (Wittmann et al. 2010). AROME-EPS produces precipitation

627    events, which are mostly too small and/or too peaked, whereas precipitation objects

628    in ALADIN-LAEF are too large and flat. This is particularly true for days with strong

629    synoptic forcing and for flat terrain. The structure score for ALADIN-LAEF further

630    shows a pronounced diurnal variation for region *West,* where precipitation events are

631    too large during the day (0900 – 1500 UTC), but more realistic during evening and

632    nighttime. In region *Northeast* and weak synoptic forcing*,* on the contrary, there is a

633    rather damped diurnal variation*.* This is a sign that precipitation events emerge too

634    early and grow too large over the mountains, whereas over flat land, they are too flat

635    and too widespread during the whole day. AROME-EPS generally shows better

636    agreement with the observed precipitation structures than ALADIN-LAEF during noon

637    (1200 - 1500 UTC) while objects are much too small during the rest of the day. Only

638    on days with strong synoptic forcing and over mountainous terrain does AROME-

639    EPS mostly underestimate the dimension of precipitation events. Over flat land,

640    structure scores are variable on a low level for AROME-EPS, but do not show a

641    perfect daily cycle.

642    In most instances, the amplitude component reflects the findings shown in Figure 5,

643    being more apparent for days with weak than for days with strong synoptic forcing.

644    For both EPS models, an overestimation occurs during noon over mountainous

645    terrain (region *West,* Figure 7), which is associated with the early onset of convection

646    for ALADIN-LAEF and with the overestimation of precipitation amounts in AROME-

647    EPS. In region *Northeast* (Figure 8), the agreement seems to be much better for

648    days with strong synoptic forcing than for weak synoptic forcing. However, amplitude

649    score measures the agreement in terms of the percentage share of precipitation

650    amounts. Hence, if the amounts are on a much lower level as in the case of weak

651    synoptic forcing, amplitude scores appear worse. The large amplitude errors in

652    Figures 8c and 8d are, therefore, more dependent on the time shift between

653    simulated and observed peaks of precipitation intensities than on the absolute

654    amount of maximum precipitation intensities, which are fairly well captured.

655    The location score in both regions provided by the SAL shows not as much variability

656    as the other two components. Nevertheless, an investigation of the distances of

657    observed and forecast centers of mass for the precipitation events can provide useful

658    information. Figures 9a and 9b show the mean distances for objects pertaining to

659    precipitation thresholds of 0.1 mm / 3 h and of 2 mm / 3 h for days with strong

660    synoptic forcing, respectively. In general, it can be stated that the distances get

661    shorter with increasing thresholds. This indicates that both ALADIN-LAEF and

662    AROME-EPS are more successful for more intense precipitation events. On the other

663    hand, precipitation objects with very low intensities can be either very small and

664    randomly distributed, which is difficult to predict, or very large, which is easier to

665    predict or detect.

666    For higher thresholds, Figure 9b shows that the distances have more variability with

667    time. Although distances are short for earlier hours of the forecast (and the first half

668    of the day), they increase for later forecast hours and reach a maximum at +21 h

669    (2100 UTC). This effect is much greater in ALADIN-LAEF than in AROME-EPS and it

31

670     is remarkable that it happens very late in the day, much later than the main peak of

671     precipitation shown in Figure 5. The reason could be that the precipitation cells are

672     captured well when they are in a mature and well developed state. Their further

673     development or collapse seems to be better simulated in AROME-EPS. This should

674     be connected to the prognostic (and explicit) treatment of the atmospheric variables

675     describing the evolution of convective activity in AROME. A convection

676     parameterization, in particular, a diagnostic convection scheme (as it is used for

677     some members of ALADIN-LAEF) has more deficiencies in simulating the life cycle of

678     convective objects properly than is the case for AROME. In addition, the non-

679     hydrostatic dynamics, higher resolution and better representation of turbulence and

680     microphysical interactions in the model physics might lead to a more realistic decay

681     of convection in AROME-EPS.

682

683     *iii.) Fractions Skill Score*

684     The fractions skill score (FSS) indicates how well the ensemble systems predict

685     precipitation at different spatial scales. The grid box widths (1 km – 21 km,

686     corresponding to areas of 1 km$^2$ – 441 km$^2$) have been selected to investigate the

687     performance of models at very fine scales, near the resolution of the analyzed

688     observations of INCA. At these scales models have difficulties to reach the level of

689     *usefulness* (i.e. the *target skill* as defined in Roberts and Lean 2008), which can be

690     expected at larger scales. Nevertheless, it is interesting to examine how FSS values

691     change with increasing precipitation thresholds.

692  Figures 10a and 10b compare the fractional skill scores for days with strong synoptic

693  forcing and days with weak forcing. FSS values are greater (~factor 2) for strong

694  synoptic forcing than for weak synoptic forcing, since for the latter, precipitation

695  events are generally less structured which lead to the lower level of skill.

696  For all weather situations, ALADIN-LAEF shows better values for the lowest

697  thresholds of 0.1 mm and 0.5 mm. The converse result is observed for higher

698  thresholds above 2 mm. For 5 mm / 3 h ALADIN-LAEF has hardly any skill on the

699  very fine scales for days with weak synoptic forcing. This means that small, scattered

700  showers and thunderstorms, which typically occur on these days, cannot be

701  simulated well by the model with coarser model resolution. In AROME-EPS there is

702  at least a certain skill for small intense precipitation events, although it is not at a

703  level considered as reliable.

704  In the previous sections, the discussion provided an overview on the whole 3 months

705  period. In the following section, evaluations focus on a single selected day. This is

706  done in order to show the forecast behavior of the ensembles in a concrete weather

707  situation exemplarily.

708  *c. Case study*

709  A typical convective day with weak synoptic forcing is selected to show the evolution

710  of precipitation in AROME-EPS and ALADIN-LAEF in more detail. Here more

711  emphasis is put on the observation of the numbers, volumes, and distribution of the

712  precipitation objects.

713    Figure 11 illustrates the precipitation at different times of 29 April 2014 of INCA

714    analyses and the ensemble means of AROME-EPS and ALADIN-LAEF. On this day,

715    continuous light rain was reported in Austria's mountainous terrain, near the main

716    Alpine ridge during the morning hours as shown in the first row of Figure 11. At the

717    same time the lowlands in the east and north were dry. In the lowlands, precipitation

718    activities in terms of small showers started from approximately 1100 UTC in second

719    row of Figure 11. Over the course of the day the focus of precipitation was

720    increasingly shifted to the flat lands in the North, East, and Southeast of Austria as

721    well as to Slovenia and Northern Italy. The peak rain intensity was around 1500 UTC,

722    shown at 1400 UTC in third row of Figure 11. Rain in the inner alpine areas had

723    diminished. In contrast, the showers in the flat regions continued until the time of

724    sunset. Then their activity also weakened, which is visible in the bottom row of Figure

725    11.

726    Figure 12 gives the characteristics of the precipitation forecasts of ALADIN-LAEF and

727    AROME-EPS, such as the temporal evolution of the mean areal precipitation in

728    Figure 12a, the number of precipitation objects in Figure 12b, and the temporal

729    evolution of the SAL scores in Figure 12c. For the selected day, precipitation

730    amounts for the region *Austria* are slightly underestimated by the both ensemble

731    systems. Further, only a minor fraction of ensemble members reach the observed

732    precipitation intensities at noon. By investigating the structures of the precipitation

733    forecasts, further insight into the behavior of the ensemble systems is provided. The

734    number and volume of precipitation objects describe how models perform in a spatial

735    context. In this respect, AROME-EPS clearly shows more ability to replicate the real

736  spatial structure of precipitation. Although the number of objects in the region *Austria*

737  is too low during the first forecast hours, the further development as observed by the

738  INCA analysis in Figure 12b is described well. In the ALADIN-LAEF forecast the

739  number of precipitation objects is very low, mostly a product of the lower resolution.

740  The volumes of the precipitation events are in direct connection with their number

741  (not shown). ALADIN-LAEF overestimates the volumes to the same degree as it

742  underestimates their numbers. However, it shows a clear diurnal variation of the

743  volumes with a maximum around noon, which is not indicated by AROME.

744  The fact that ALADIN-LAEF tends to produce fewer but larger precipitation objects

745  does not lead to worse verification statistics for ALADIN-LAEF. On the contrary, in

746  most regions the hit rate is higher for ALADIN-LAEF than for AROME-EPS and the

747  number of missed events is lower. AROME-EPS, on the other hand outperforms

748  ALADIN-LAEF in terms of correct negatives and false alarms (not shown).

749  These results are also reflected in the temporal evolution of SAL-scores in Figure

750  12c. As expected, the structure score S is too high for ALADIN-LAEF, due to the

751  overestimation of the volumes of precipitation objects. At the same time, however,

752  AROME-EPS produces a low S score which means that it still produces too small

753  and peaked precipitation objects compared to INCA.

754  Interestingly, there is a late peak in the S score between 26-28 hours lead time in

755  both models, which follows a short minimum at 25 hours lead time. This is also

756  slightly reflected in the A score. The sequence of minimum and peak is related to a

757  nightly shower, which was also simulated by the ensembles, but with a delay of

758  approximately 2 hours. The location or L-score is rather constant in time for both

759   ensemble models. This means that they were able to reproduce the changing spatial

760   focus and distribution of precipitation during the day.

761   **5. Summary and conclusions**

762   In this paper we investigate the forecast performance of the 2.5 km convection-

763   permitting ensemble AROME-EPS by comparison with the regional 11 km ensemble

764   ALADIN-LAEF to reveal the benefit provided by a CP EPS. The regional EPS,

765   ALADIN-LAEF, involves several sources of forecast perturbations, such as initial

766   condition perturbations by blending ECMWF-EPS with ALADIN-LAEF breeding

767   vectors and assimilation of perturbed surface observations, and a multi-physics

768   scheme. The high-resolution, convection-permitting AROME-EPS solely performs

769   downscaling of the ALADIN-LAEF forecasts. The performance of the ensembles is

770   evaluated for a 3-month period during the convective season of 2011 and for a

771   typical convective day in April 2014 with a special focus on precipitation events in

772   mountainous terrain and lowland regions. The aim is to show whether the

773   convection-permitting ensemble provides benefits to the regional ensemble with deep

774   convection parameterization. The evaluation is conducted using a combination of

775   standard deterministic and probabilistic verification scores and selected spatial

776   verification measures. The former are applied on several main forecast parameters

777   for surface and upper levels, the latter – according to their definition – only for

778   precipitation.

779   The forecast quality for the main meteorological parameters (except precipitation) for

780   the surface and selected upper levels is strongly dependent on the model bias and is

781   rather balanced, except for diurnal variations near the surface. However,

782    characteristic differences are revealed by the investigation of the precipitation

783    forecasts. A known drawback of models using deep convection schemes proves true,

784    which is the premature onset of precipitation in the daily cycle by ALADIN-LAEF (see

785    e.g. Wittmann et al., 2010; Weusthoff et al., 2010). On the other hand, an

786    overestimation of precipitation intensities at the peak of convection activities by

787    AROME-EPS is also confirmed, which has been assumed in previous validations.

788    Both of these properties are found to be more pronounced in mountainous than in flat

789    regions.

790    ALADIN-LAEF shows skill in the prediction of probabilities for low precipitation

791    thresholds, i.e. to distinguish between *rain* and *no rain*. This is also true for small

792    scales, but it is again dependent on the time of day, as the early onset of precipitation

793    has a negative influence on the verification scores. AROME-EPS, on the other hand,

794    has a better ability to capture the diurnal cycle of convective precipitation, especially

795    over mountainous terrain. At small spatial scales, it further demonstrates better

796    performance for higher precipitation thresholds. The results of the evaluations in this

797    study lead to the conclusion, that the convection permitting ensemble is more skillful

798    on the precipitation forecast than its mesoscale counterpart, the regional ensemble.

799    The positive impact is larger for the mountainous areas than for the lowlands.

800    Nevertheless, the knowledge of which precipitation situations can be better modeled

801    by the convection-permitting ensemble is important to have. For many applications,

802    e.g. for large-scale extreme events, such as the Central Europe flooding event of

803    2013, the best solution will be a combination of both systems: the coarser ensembles

804    with longer forecast range for (pre)-warnings, and the convection-permitting

805 ensemble for the detailed specification of the expected event. Regarding different

806 time and length-scales in that way could lead to the generation of *seamless* forecast

807 products (e.g. Drobinski et al. 2014, Vitart et al. 2008).

808 This study is considered as initial point for further investigations and improvement of

809 the convection-permitting ensemble AROME-EPS. The low spread of the prevailing

810 AROME-EPS version is a clear drawback compared to ALADIN-LAEF. Therefore,

811 future enhancements of AROME-EPS will involve components, which will

812 presumably increase ensemble spread. Among those upgrades will be ensemble

813 data assimilation and physics perturbations (multi-model and stochastic). The

814 expectation with these components is that forecast errors will be reduced, and that a

815 more realistic simulation of forecast uncertainties will be achieved.

816 **6. Code and/or data availability**

817 The ALADIN-LAEF and AROME codes including all related intellectual property

818 rights, are owned by the members of the LACE consortium and ALADIN consortium.

819 Access to the ALADIN-LAEF and AROME systems, or elements thereof, can be

820 granted upon request and for research purposes only. INCA and INCA data are only

821 available subject to a licence agreement with ZAMG.

822

823 Acknowledgments

824 We gratefully acknowledge all the LACE/ALADIN/HIRLAM colleagues who have

825 contributed to the development of AROME. ECMWF has provided the computer

826 facilities and technical help implementing ALADIN-LAEF and AROME-EPS on the

827 ECMWF HPCF.

828 REFERENCES

829 Ahijevych D., E. Gilleland, B. Brown, and E. Ebert, 2009: Application of spatial

830 forecast verification methods to gridded precipitation forecasts. *Wea. Forecasting,*

831 **24**, 1485–1497.

832 Aligo A. E., W. A. Gallus Jr., and M. Segal, 2009: On the Impact of WRF Model

833 Vertical Grid Resolution on Midwest Summer Rainfall Forecasts. *Wea. Forecasting*,

834 **24**, 575-594.

835 Barthlott C., R. Burton, D. Kirshbaum, K. Hanley, R. Richard, J. P. Chaboreau, J.

836 Trentmann, B. Kern, H.-S. Bauer, T. Schwitalla, C. Keil, Y. Seity, A. Gadian, A. M.

837 Blyth, S. Mobbs, C. Flamant, and J. Handwerker, 2011:Initiation of deep convection

838 at marginal instability in an ensemble of mesoscale models: A case-study from

839 COPS. *Quart. J. Roy. Meteor. Soc.,* **137**, 118–136.

840 Bauer H.S., T. Weusthoff, M. Dorninger, V. Wulfmeyer, T. Schwitalla, T. Gorgas, M.

841 Arpagaus, and K. Warrach-Sagi, 2011: Predictive skill of a subset of models

842 participating in D-PHASE in the COPS region. *Q. J. R. Meteorol. Soc*. **137**, 287-305.

843 Bechtold, P., N. Semane, P. Lopez, and J.-P. Chaboureau, A. Beljaars, N. Bormann,

844 2013: Breakthrough in forecasting equilibrium and non.equilibrium convection.

845 *ECMWF Newsletter,* **136**, 15-22.

846 Bowler, N. E., A. Arribas, K. R. Mylne, K. B. Robertson, and S. E. Beare, 2008: The

847    MOGREPS short-range ensemble prediction system. *Quart. J. Roy. Meteor. Soc.*,

848    **134**, 703–722.

849    Ben Bouallégue, Z., S. E. Theis, and C. Gebhardt, 2013: Enhancing COSMO-DE

850    ensemble forecasts by inexpensive techniques. *Meteor. Z.*, **22**, 49–59.

851    Bougeault, P., 1985: A simple parameterization of the large-scale effects of cumulus

852    convection. *Mon. Wea. Rev.,* **113**, 2108–2121.

853    Bouttier, F., B. Vié, O. Nuissier, and L. Raynaud, 2012: Impact of Stochastic Physics

854    in a Convection-Permitting Ensemble. *Mon. Wea. Rev.*, **140**, 3706-3721.

855    Brier, G. W., 1950: Verification of forecasts expressed in terms of probability. *Mon.*

856    *Wea. Rev.*, **78**, 1–3.

857    Bryan, G. H., J. C. Wyngaard, and J. M. Fritsch, 2003: Resolution requirements for

858    the simulation of deep moist convection. *Mon. Wea. Rev.*, **131,** 2394–2416.

859    Bukovsky, M. S., J. S. Kain and M. E. Baldwin, 2006: Bowing convective systems in

860    a popular operational model: Are they for real? *Wea. Forecasting*, **21,** 307–324.

861    Caron, J., 2013: Mismatching perturbations at the lateral boundaries in limited-

862    areaensemble forecasting: A case study. *Mon. Wea. Rev.,* **141**, 356–374.Casati B. L,

863    L. J. Wilson, D. B. Stephenson, P. Nurmi, A. Ghelli, M. Pocernich, U. Damrath, E. E.

864    Ebert, B. G. Brown, and S. Mason, 2008: Review forecast verification: current status

865    and future directions. *Meteor. Appl.*, **15**, 3–18.

866    Clark, A. J., W. A. Gallus Jr., and T.-C. Chen, 2007: Comparison of the Diurnal

867    Precipitation Cycle in Convection-Resolving and Non-Convection-Resolving

868    Mesoscale Models. *Mon. Wea. Rev.,* **135**, 3456-3473.

869    Clark, A. J., W. A. Gallus Jr., M. Xue, and F. Kong, 2009: A Comparison of

870    Precipitation Forecast Skill between Small Convection-Allowing and Large

871    Convection-Parameterizing Ensembles. *Wea. Forecasting,* **24**, 1121-1140.

872    Clark, A. J., J. S. Kain, D. J. Stensrud, M. Xue, F. Kong, M. C. Coniglio, K. W.

873    Thomas, Y. Wang, K. Brewster, J. Gao, X. Wang, S. J. Weiss and J. Du, 2011:

874    Probabilistic Precipitation Forecast Skill as a Function of Ensemble Size and Spatial

875    Scale in a Convection-Allowing Ensemble. *Mon. Wea. Rev.,* **139**: 1410-1418.

876    Davis, C. A., K. W. Manning, R. E. Carbone, S. B. Trier, and J. D. Tuttle, 2003:

877    Coherence of warm season continental rainfall in numerical weather prediction

878    models. *Mon. Wea. Rev.*, **131,** 2667–2679.

879    Davison, A.C. and D.V. Hinkley, 1997: Bootstrap Methods and their applications –

880    Cambridge University Press, Cambridge, UK, 193 f.

881    Done, J. M., G. C. Craig, S. L. Gray, P. A. Clark, and  M. E. B. Gray, 2006:

882    Mesoscale simulations of organized convection: Importance of convective

883    equilibrium. *Quart. J. Roy. Meteor. Soc.,* **132**, 737–756.

884    Drobinski, P., and Coauthors, 2014: HyMeX: A 10-Year Multidisciplinary Program on

885    the Mediterranean Water Cycle*. Bull. Amer. Meteor. Soc.*, **95**, 1063-1082.

886  Ferro, C.A.T., 2007: A probability model for verifying deterministic forecasts of
887  extreme events. *Wea. Forecasting*, **22**, 1089–1100.

888  Gebhardt, C., S. E. Theis, M. Paulat and Z. Ben Bouallègue, 2011: Uncertainties in
889  COSMO-DE precipitation forecasts introduced by model perturbations and variation
890  of lateral boundaries. *Atmos. Res.*, **100**, 168-177

891  Geleyn, J.-F., B. Catry, Y. Bouteloup, and R. Brožková, 2008: A statistical approach
892  for sedimentation inside a microphysical precipitation scheme. *Tellus*, **60A**, 649–662,
893  doi:10.1111/j.1600-0870.2008.00323.x.

894  Gerard, L., J.-M. Piriou, R. Brožkova, J.-F. Geleyn, and D. Banciu, 2009: Cloud and
895  precipitation parameterization in a meso-gamma scale operational weather prediction
896  model*. Mon. Wea. Rev.*, **137**, 3960–3977.

897  Gilleland, E., D. A. Ahijevych, B. G. Brown, and E. E. Ebert, 2010: Verifying forecasts
898  spatially. *Bull. Amer. Meteor. Soc.*, **47**, 1365–1373.

899  Gneiting, T., and A. E. Raftery, 2007: Strictly Proper Scoring Rules, Prediction and
900  Estimation. *Journal of the American Statistical Association,* **102,** 359-378.

901  Haiden, T., and G. Pistotnik, 2009: Intensity-dependent parameterization of elevation
902  effects in precipitation analysis. *Adv. Geosci.*, **20**, 33-38.

903  Haiden, T., A. Kann, C. Wittmann, G. Pistotnik, B. Bica, and C. Gruber, 2011: The
904  Integrated Nowcasting through Comprehensive Analysis (INCA) System and Its
905  Validation over the Eastern Alpine Region. *Wea. Forecasting*, **26**, 166-183.

906 Haiden, T., L. Magnusson, I. Tsonevsky, F. Wetterhall, L. Alfieri, F. Pappenberger, P.

907 de Rosnay, J. Muñoz-Sabater, G. Balsamo, C. Albergel, R. Forbes, T. Hewson, S.

908 Malardel, and D. Richardson, 2014: ECMWF forecast performance during the June

909 2013 flood in Central Europe. *ECMWF – Technical Memorandum,* **723**,

910 http://old.ecmwf.int/publications/library/ecpublications/pdf/tm/701-800/tm723.pdf (Sep

911 2, 2014).

912 Hanley, K. E., D. J. Kirshbaum, N. M. Roberts and G. Leoncini, 2013: Sensitivities of

913 a Squall Line over Central Europe in a Convective-Scale Ensemble. *Mon. Wea. Rev.*,

914 **141**, 112-133.

915 Hersbach, H., 2000: Decomposition of the Continuous Ranked Probability Score for

916 Ensemble Prediction Systems. *Wea. Forecasting,* **15,** 559-570.

917 Jolliffe, I., 2007: Uncertainty and inference for verification measures. *Wea.*

918 *Forecasting,* **22**, 637–650.

919 Kühnlein, C., C. Keil, G. C. Craig, and C. Gebhardt, 2014: The impact of downscaled

920 initial condition perturbations on convective-scale ensemble forecasts of precipitation.

921 *Quart. J. Roy. Meteor. Soc.,* **140**, 1552–1562.

922 Liu, C., M. W. Moncrieff, J. D. Tuttle, and R. E. Carbone, 2006 : Explicit an

923 Parameterized Episodes of Warm-Season Precipitation over the Continental United

924 States. *Adv. Atmos. Sci.*, **23**, 91-105.

925 Masson, V., 2000: A physically-based scheme for the urban energy budget in

926 atmospheric models. *Bound.-Layer Meteor.*, **94**, 357–397.

927   Mlawer, E. J., S. J. Taubman, P. D. Brown, M. J. Iacono, and S. A. CLough, 1997 :

928   Radiative transfer for inhomogeneous atmospheres: RRTM, a validated correlated-k

929   model for the longwave. *J. Geophys. Res.*, **102D**, 16663–16682.

930   Morcrette, J.-J., 1991 : Radiation and cloud radiative properties in the ECMWF

931   operational weather forecast model. *J. Geophys. Res.*, **96D**, 9121–9132.

932   Murphy, A. H., 1973: A new vector partition of the probability score. *J. Appl. Meteor.*,

933   **12**, 595–600.

934   Nurmi, P., 2003: Recommendations on the verification of local weather forecasts.

935   *ECMWF Technical Memoranda,* **430**, 19 pp. [available online at

936   http://old.ecmwf.int/publications/library/do/references/show?id=86094] (Oct 20,

937   2014).

938   Peralta, C., Z. Ben Bouallègue, S. E. Theis, C. Gebhardt, and M. Buchhold, 2012:

939   Accounting for initial condition uncertainties in COSMO-DE-EPS. *J. Geophys. Res.,*

940   **117**, 1–13, doi: 10.1029/2011JD016581.

941   Pergaud, J., V. Masson, V., and S. Malardel, 2009: A parameterization of dry

942   thermals and shallow cumuli for mesoscale numerical weather prediction, *Bound.-*

943   *Layer Meteor.*, **132**, 83–106.

944   Pinty, J. P., and Jabouille, P., 1998: A mixed phase cloud parameterization for use in

945   a mesoscale nonhydrostatic model: Simulations of a squall line and of orographic

946   precipitation. *Preprints, Conf. on Cloud Physics*, Everett, WA, Amer. Meteor. Soc.,

947   217–220.

948    Richard E., Buzzi A., and Zängl G., 2007. Quantitative precipitation forecasting in the

949    Alps: The advances achieved by the Mesoscale Alpine Programme. *Q. J. R.*

950    *Meteorol. Soc.* **133**: 831–846.

951    Ritter, B., and J.-F. Geleyn, 1992: A comprehensive radiation scheme for numerical

952    weather prediction models with potential applications in climate simulations. *Mon.*

953    *Wea. Rev.*, **120**, 303–325.

954    Roberts, N. M., and H. W. Lean, 2008: Scale-selective verification of rainfall

955    accumulations from high-resolution forecasts of convective events. *Mon. Wea. Rev.*,

956    **136**, 78–97.

957    Romine, G. S., C. S. Schwartz, J. Berner, K. R. Fossell, C. Snyder, J. L. Anderson,

958    M. L. Weisman, 2014: Representing Forecast Error in a Convection-Permitting

959    Ensemble System. *Mon. Wea. Rev.,* **142**, 4519-4541.

960    Rotach, M. W., and Coauthors including and T. Gorgas and Y. Wang, 2009: MAP D-

961    PHASE real time demonstration of weather forecast quality in the Alpine region. *Bull.*

962    *Amer. Meteor. Soc.*, **90**: 1321-1336.

963    Schumacher, R. S., A. J. Clark, M. Xue, and F. Kong, 2013: Factors Influencing the

964    Development and Maintenance of Nocturnal Heavy-Rain-Producing Convective

965    Systems in a Storm-Scale Ensemble. *Mon. Wea. Rev.*, **141**: 2778-2801.

966    Schumacher, R. S., and A. J. Clark, 2014: Evaluation of ensemble configurations for

967    the analysis and prediction of heavy-rain-producing mesoscale convective systems.

968    *Mon. Wea. Rev.*, **e-View**, doi: http://dx.doi.org/10.1175/MWR-D-13-00357.1.

969  Schwartz C. S, J. S. Kain, S. J. Weiss, M. Xue, D. R. Bright, F. Kong, K. W. Thomas,

970  J. J. Levit, M. C. Coniglio, M. S. Wandishin, 2010: Toward Improved Convection-

971  Allowing Ensembles: Model Physics Sensitivities and Optimizing Probabilistic

972  Guidance with Small Ensemble Membership. *Wea. Forecasting*, **25**, 263–280.

973  DOI:10.1175/2009WAF2222267.1.

974  Schwartz, C. S., G. S. Romine, K. R. Smith, and M. L. Weisman, 2014:

975  Characterizing and optimizing precipitation forecasts from a convection-permitting

976  ensemble initialized by a mesoscale ensemble kalman filter. *Wea. Forecasting*, **29**,

977  1295–1318.

978  Seity, Y., P. Brousseau, S. Malardel, G. Hello, P. Bénard, F. Bouttier, C. Lac, and V.

979  Masson, 2011: The AROME-France Convective-Scale Operational Model. *Mon.*

980  *Wea. Rev.* **139**: 976–991.

981  Taillerfer, F., 2002: CANARI – Technical Documentation - Based on ARPEGE cycle

982  CY25T1      (AL25T1      for      ALADIN),      [available      online      at:

983  http://www.cnrm.meteo.fr/gmapdoc/IMG/ps/canari_doc_cy25t1.ps   (cited   Dec   14,

984  2015)]

985  Taraphdar, S., P. Mukhopadhyay, L. R. Leung, F. Zhang, S. Abhilash, and B. N.

986  Goswami, 2014: The role of moist processes in the intrinsic predictability of Indian

987  Ocean   cyclones,   *J.   Geophys.   Res.   Atmos.*,   **119**,   8032-8048,

988  doi:10.1002/2013JD021265

989  Tennant, W., 2015: Improving initial condition perturbations for MOGREPS-UK.

990  *Quart. J. Roy. Meteor. Soc.,* DOI: 10.1002/qj.2524. Online publication date: 1-Feb-2015.

991   Theis, S. E., A. Hense, U. Damrath, 2005: Probabilistic precipitation forecasts from a

992   deterministic model: a pragmatic approach. *Meteor. Appl.* **12**, 257–268.

993   DOI:10.1017/S1350482705001763.

994   Toth, Z., and E. Kalnay, 1993: Ensemble forecasting at NMC: the generation of

995   perturbation. *Bull. Amer. Meteor. Soc.*, **74**, 2317-2330.

996   UK Met Office, July 2014 (cited Oct 21, 2014): Benefits of high resolution ensemble

997   forecasts. [available online at http://www.metoffice.gov.uk/research/news/2014/high-

998   resolution-ensembles]

999   Vana F., P. Benard, J.-F. Geleyn, A. Simon, and Y. Seity, 2008: Semi-Lagrangian

1000   advection scheme with controlled damping: an alternative to nonlinear horizontal

1001   diffusion in a numerical weather prediction model. *Quart. J. Roy. Meteor. Soc.,* **134**,

1002   523–537.

1003   Vié, B., G. Molinié, O. Nussier, B. Vincendon, V. Ducrocq, F. Bouttier, and E.

1004   Richard, 2012: Hydro-meteorological evaluation of a convection-permitting ensemble

1005   prediction system for Mediterranean heavy precipitating events. *Nat. Hazards Earth*

1006   *Syst. Sci.*, **12**: 2631–2645.

1007

1008   Vitart F., R. Buizza, M. A. Balmaseda, G. Balsamo, J.-R. Bidlot, A. Bonet, M.

1009   Fuentes, A. Hofstadler, F. Molteni, and T. N. Palmer, 2008: The new VarEPS–

1010   monthly forecasting system: A first step towards seamless prediction. *Quart. J. Roy.*

1011   *Meteor. Soc.*, **134**, 1789–1799.

1012 Wang, Y., A. Kann, M. Bellus, J. Pailleux, and C. Wittmann, 2010: A strategy for
1013 perturbing surface initial conditions in LAMEPS. *Atmos. Sci. Let.,* **11**, 108-113.

1014 Wang, Y., M. Bellus, C. Wittmann, M. Steinheimer, F. Weidle, A. Kann, S. Ivatek-
1015 Šahdan, W. Tian, X. Ma, S. Tascu, and E. Bazile, 2011: The Central European
1016 limited-area ensemble forecasting system: ALADIN-LAEF. *Quart. J. Roy. Meteor.*
1017 *Soc.,* **137**, 483–502.

1018 Wang, Y., S. Tascu, F. Weidle, and K. Schmeisser, 2012: Evaluation of the Added
1019 Value of Regional Ensemble Forecasts on Global Ensemble Forecasts. Wea.
1020 Forecasting, **27**, 972-987.

1021 Wang Y., M. Bellus, J.-F. Geleyn, X. Ma, W. Tian, and F. Weidle, 2014: A New
1022 Method for Generating Initial Condition Perturbations in a Regional Ensemble
1023 Prediction System: Blending*. Mon. Wea. Rev.*, **142**, 2043-2059.

1024 Weckwerth, T., L. Bennett, L. Miller, J. Van Baelen, P. Di Girolamo, A. Blyth, and T.
1025 Hertneky, 2014: An Observational and Modeling Study of the Processes Leading to
1026 Deep, Moist Convection in Complex Terrain. *Mon. Wea. Rev.*, **142**, 2687-2708.

1027 Weidle, F., Y. Wang, W. Tian and T. Wang, 2013: Validation of Strategies using
1028 Clustering Analysis of ECMWF EPS for Initial Perturbations in a Limited Area Model
1029 Ensemble Prediction System. *Atmosphere-Ocean*, **51**, 284-295.

1030 Weisman, M. L., W. C. Skamarock, and J. B. Klemp, 1997: The resolution
1031 dependence of explicitly modeled convective systems. *Mon. Wea. Rev.*, **125,** 527–
1032 548.

1033  Wernli, H., M. Paulat, M. Hagen, and C. Frei, 2008: SAL - A Novel Quality Measure

1034  for the Verification of Quantitative Precipitation Forecasts. *Mon. Wea. Rev.,* **136**,

1035  4470–4487.

1036  Weusthoff, T., F. Ament, M. Arpagaus, and M. W. Rotach, 2010: Assessing the

1037  Benefits of Convection-Permitting Models by Neighborhood Verification: Examples

1038  from MAP D-PHASE. *Mon. Wea. Rev.*, **138**, 3418–3433.

1039  Wilks, D. S., 1997: Resampling hypothesis testing for autocorrelated fields.

1040  *J.Climate,* 10, 65-82.

1041  Wilks, D. S., 2006: *Statistical Methods in the Atmospheric Sciences*, 2nd Ed.,

1042  London, Academic Press, 627 pp.

1043  Wittmann, C., T. Haiden, and A. Kann, 2010: Evaluating multi-scale precipitation

1044  forecasts using high resolution analysis. *Adv. Sci. Res.*, **4**, 89-98, DOI:10.5194/asr-4-

1045  89-2010.

1046  Wulfmeyer V., and Coauthors, 2008: The Convective and Orographically induced

1047  Precipitation Study: A research and development project of the World Weather

1048  Research Program for improving quantitative precipitation forecasting in low-

1049  mountain regions. *Bull. Am. Meteorol. Soc.* **89**: 1477–1486,

1050  DOI:10.1175/2008BAMS2367.1.

1051  Wulfmeyer, V., and Coauthors including T. Gorgas and Y. Wang, 2011: The

1052  Convective and Orographically-induced Precipitation Study (COPS): the scientific

1053    strategy, the field phase, and research highlights. *Quart. J. Roy. Meteor. Soc.* **137**:

1054    3–30.

1055    Xue, M., and Coauthors, 2007: CAPS realtime storm-scale ensemble and high-

1056    resolution forecasts as part of the NOAA Hazardous Weather Testbed 2007 spring

1057    experiment. *Extended Abstracts, 22nd Conference on Weather Analysis and*

1058    *Forecasting/18th Conference on Numerical Weather Prediction*, Park City, UT. Amer.

1059    Meteor. Soc.,[Available onlineat http://ams.confex.com/ams/pdfpapers/124587.pdf].

1060    Xue, M., and Coauthors, 2009: CAPS realtime multi-model convection-allowing

1061    ensemble and 1-km convection-resolving forecasts for the NOAA Hazardous

1062    Weather Testbed 2009 spring experiment. Preprints, *23rd Conf. on Weather Analysis*

1063    *and Forecasting/19$^{th}$ Conf. on Numerical Weather Prediction,* Omaha, NE, Amer.

1064    Meteor. Soc., 16A.2. [Available online at http://ams.confex.com/ams/ 23WAF19NWP/

1065    techprogram/paper_154323.htm.]

1066    Zhu, Y., Z. Toth, R. Wobus, D. Richardson, and K. Mylne, 2002: The Economic Value

1067    Of Ensemble-Based Weather Forecasts. *Bull. Amer. Meteor. Soc.* **83**, 73-83.

1068

1069

1070

|  | ALADIN-LAEF | AROME-EPS |
|---|---|---|
| Ensemble size | 16+1 members | 16 members |
| Horizontal resolution | 11 km | 2.5 km |
| Vertical resolution | 45 layers | 60 layers |
| Model time step | 450 s | 60 s |
| Coupling-Model | ECMWF-EPS | ALADIN-LAEF |
| Coupling-Update | 6 h | 3 h |
| No. of grid points | 206 x 164 | 432 x 320 |
| Forecast range | 72 h | 30 h |
| Runs/Day | 2 (0000, 1200 UTC) | 1 (0000 UTC) |

1071

1072 **Table 1**: Main characteristics of the ALADIN-LAEF and AROME-EPS.

1073

1074

1075

1076

1077

1078

**Formatiert:** Links, Abstand Nach: 0 Pt., Tabstopps: Nicht an 14,92 cm + 16 cm

1079    Figure 1: Geographic domains and topographies of a) ALADIN-LAEF, where the red

1080    frame is the output domain used for the present study, and b) AROME-EPS, which is

1081    shown by the blue frame in (a).

1082

1083    Figure 2: Locations of meteorological surface observation stations within the

1084    evaluation domain.

1085

1086    Figure 3: INCA domain and topography with the sub-domains, which are used for the

1087    evaluation.

1088

1089    Figure 4: Bias (left panel) and CRPS (right panel) for 2m relative humidity (top), 2m

1090    temperature (middle) and 10m wind speed (bottom) for the period of May 15 –

1091    August 15, 2011 of AROME-EPS (dotted line) and ALADIN-LAEF (solid line), both

1092    verified over the AROME-domain. Lead times, which are marked with asterisks (*)

1093    indicate results with significant differences between the ensembles.

1094    Figure 4: Bias (left panel) and CRPS (right panel) for 2m relative humidity (top), 2m

1095    temperature (middle) and 10m wind speed (bottom) for the period of May 15 –

1096    August 15, 2011 in the AROME-domain of AROME-EPS (dotted line) and ALADIN-

1097    LAEF (solid line). Lead times, which are marked with asterisks (*) indicate results

1098    with significant differences between the ensembles.

1099

1100 Figure 5: Time evolution of 3-hourly accumulated precipitation forecast for INCA

1101 (solid line), ALADIN-LAEF ensemble mean (dashed line) and AROME-EPS

1102 ensemble mean (dotted line) for regions *Austria* (top)*, West* (middle) and *Northeast*

1103 (bottom). Left panels show results for the days with strong synoptic forcing, right

1104 panels for weak synoptic forcing. The shaded areas denote the range of individual

1105 ensemble member forecasts for ALADIN-LAEF (dark grey) and AROME-EPS (light

1106 grey) respectively.

1107

1108 Figure 6: Time evolution of the Brier Score with confidence intervals (shades) for

1109 region *Austria*, AROME-EPS (dotted line) and ALADIN-LAEF (dashed line). a) strong

1110 synoptic forcing and precipitation threshold 0.1 mm / 3 h, b) weak synoptic forcing

1111 and 0.1 mm / 3 h, c) strong synoptic forcing and 2 mm / 3 h, and d) weak synoptic

1112 forcing and 2 mm / 3 h.

1113

1114 Figure 7: Time evolution of SAL scores for AROME-EPS (left) and ALADIN-LAEF

1115 (right) for different forecast ranges in region *West*. Upper panels a) and b) show

1116 results for days with strong synoptic forcing; lower panels c) and d) for weak synoptic

1117 forcing. The boxes are created based on the scores of all individual ensemble

1118 members.

1119

1120 Figure 8: Same as in Figure 7, but for region *Northeast*.

1121

1122    Figure 9: Distances [km] between the centers of mass of the precipitation objects in

1123    the forecast and analysis fields for AROME-EPS (dotted) and ALADIN-LAEF

1124    (dashed) for thresholds of a) 0.1 mm / 3 h, and b) 2 mm / 3 h.

1125

1126    Figure 10: Fractional skill scores for a) strong synoptic forcing, and b) weak synoptic

1127    forcing of AROME-EPS (dashed) and ALADIN-LAEF (solid line) for the region

1128    *Austria*. Numbers denote the precipitation thresholds [mm]. The values represent

1129    averages for all hours of lead-time.

1130

1131    Figure 11: Observed (INCA, first column) and forecast (AROME-EPS and ALADIN-

1132    LAEF, second and third column, respectively) development of precipitation on 29

1133    April 2014 shown for selected times (rows). The panels show 1-hourly accumulated

1134    precipitation sums [mm].

1135

1136    Figure 12: Characteristics of the precipitation forecasts of ALADIN-LAEF and

1137    AROME-EPS on 29 April 2014. a) Temporal evolution of the mean areal precipitation

1138    compared with INCA, and b) temporal evolution of the number of precipitation

1139    objects. Dashed and dotted lines represent the ensemble mean and grey shades the

1140    ensemble spread. c) Temporal evolution of S (structure), A (amplitude) and L

1141    (location) scores of the ensemble means of ALADIN-LAEF (black) and AROME-EPS

1142    (grey).

1143

1144

1145

1146

1147

1148

1149

1150

a)                                                    b)

1151

1152 Figure 1: Geographic domains and topographies of a) ALADIN-LAEF, where the red

1153 frame is the output domain used for the present study, and b) AROME-EPS, which is

1154 shown by the blue frame in (a).

1155

1156

1157

1158

Figure 2: Locations of meteorological surface observation stations within the evaluation domain.

1168

1169



1170

1171 Figure 3: INCA domain and topography with the sub-domains, which are used for the

1172 evaluation.

1173

1174

1175

1176

1177

1178

1179



1180

1181  Figure 4: Bias (left panel) and CRPS (right panel) for 2m relative humidity (top), 2m

1182  temperature (middle) and 10m wind speed (bottom) for the period of May 15 –

1183  August 15, 2011 ~~in the AROME-domain~~ of AROME-EPS (dotted line) and ALADIN-

1184  LAEF (solid line), both verified over the AROME-domain.~~.~~ Lead times, which are

1185  marked with asterisks (*) indicate results with significant differences between the

1186  ensembles.

1187

1188

1189



Figure 5: Time evolution of 3-hourly accumulated precipitation forecast for INCA (solid line), ALADIN-LAEF ensemble mean (dashed line) and AROME-EPS

1193   ensemble mean (dotted line) for  regions *Austria* (top)*, West* (middle) and *Northeast*

1194   (bottom). Left panels show results for the days with strong synoptic forcing, right

1195   panels for weak synoptic forcing. The shaded areas denote the range of individual

1196   ensemble member forecasts for ALADIN-LAEF (dark grey) and AROME-EPS (light

1197   grey) respectively.

1198

1199

1200

1201

1202

1203

1204



Figure 6: Time evolution of the Brier Score with confidence intervals (shades) for region *Austria*, AROME-EPS (dotted line) and ALADIN-LAEF (dashed line). a) strong synoptic forcing and precipitation threshold 0.1 mm / 3 h, b) weak synoptic forcing and 0.1 mm / 3 h, c) strong synoptic forcing and 2 mm / 3 h, and d) weak synoptic forcing and 2 mm / 3 h.

1211

1212

On the forecasting skills of a convection permitting ensemble, SCHELLANDER-GORGAS ET AL.

1213



AROME-EPS          ALADIN-LAEF

1214

1215 Figure 7: Time evolution of SAL scores for AROME-EPS (left) and ALADIN-LAEF

1216 (right) for different forecast ranges in region *West*. Upper panels a) and b) show

1217 results for days with strong synoptic forcing; lower panels c) and d) for weak synoptic

1218 forcing. The boxes are created based on the scores of all individual ensemble

1219 members.

1220

1221

Formatiert: Schriftart: 12 Pt.

1222

1223    Figure 8: Same as in Figure 7, but for region *Northeast*.

1224

1225

63

1226

1227

a)

b)

1228 Figure 9: Distances [km] between the centers of mass of the precipitation objects in

1229 the forecast and analysis fields for AROME-EPS (dotted) and ALADIN-LAEF

1230 (dashed) for thresholds of a) 0.1 mm / 3 h, and b) 2 mm / 3 h.
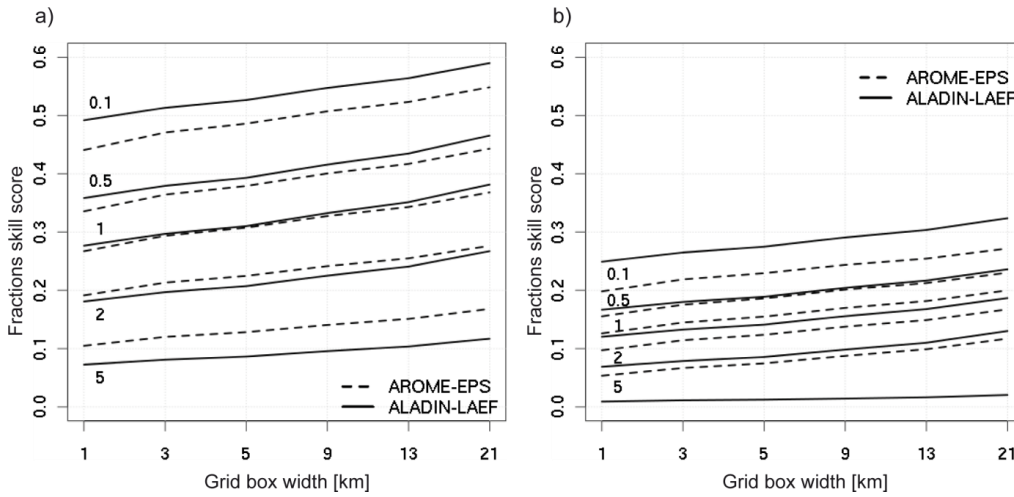
1231

1232

1233

1234

1235

1236

1237

1238

64

1239



1240

1241  Figure 10: Fractional skill scores for a) strong synoptic forcing, and b) weak synoptic

1242  forcing of AROME-EPS (dashed) and ALADIN-LAEF (solid line) for the region

1243  *Austria*. Numbers denote the precipitation thresholds [mm]. The values represent

1244  averages for all hours of lead-time.
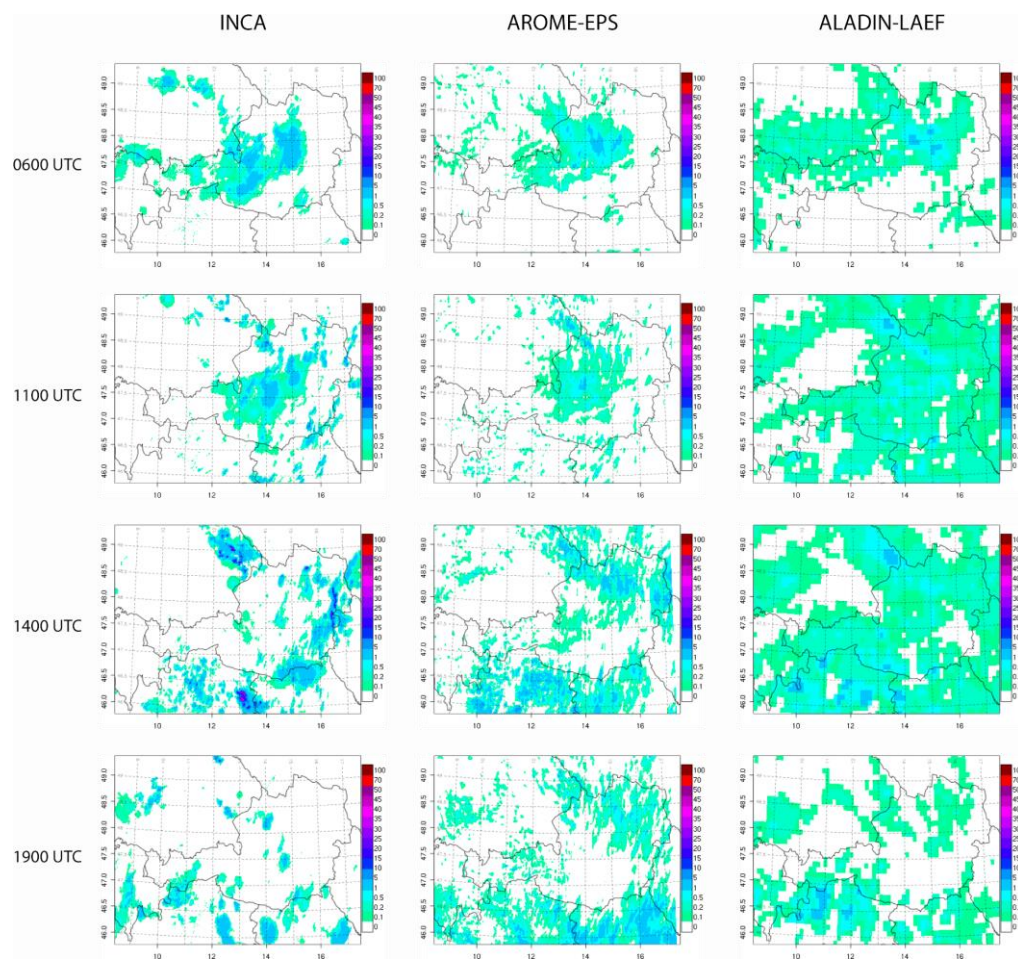
1245

1246

1247

1248

1249
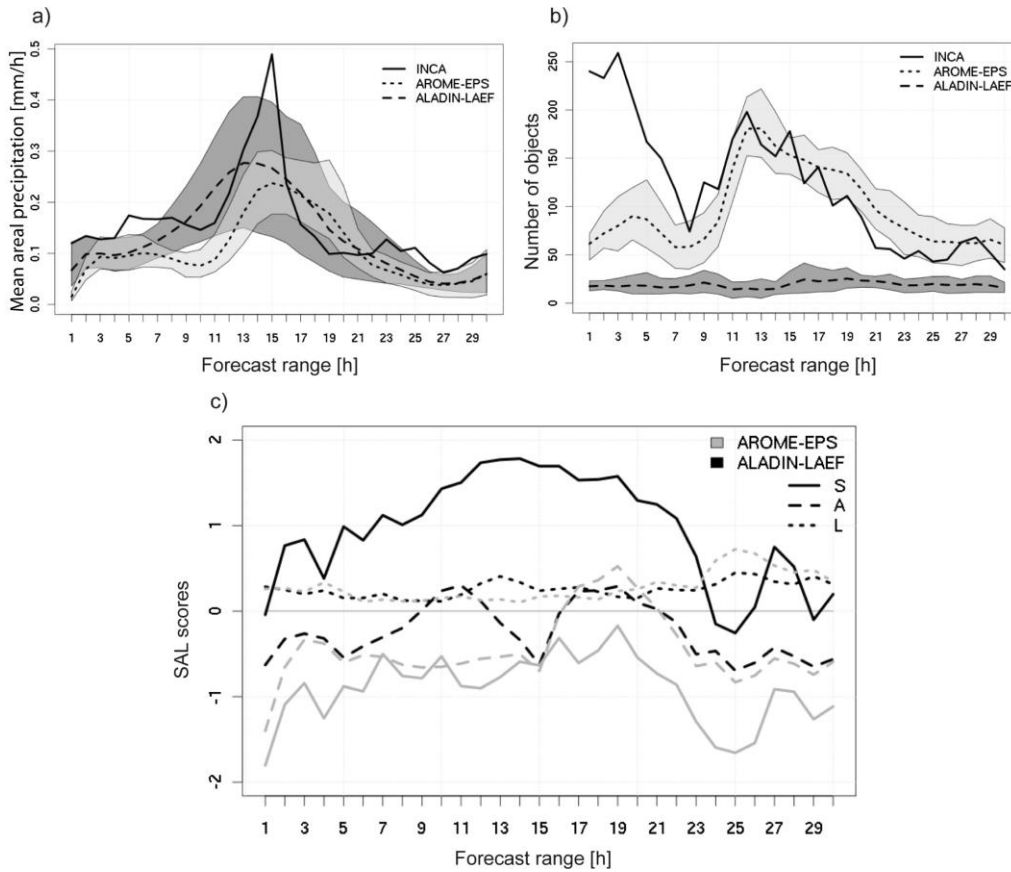
1250

1251

1252

1253

1254



1255

Figure 11: Observed (INCA, first column) and forecast (AROME-EPS and ALADIN-LAEF, second and third column, respectively) development of precipitation on 29 April 2014 shown for selected times (rows). The panels show 1-hourly accumulated precipitation sums [mm].

1260

1261



1262

Figure 12: Characteristics of the precipitation forecasts of ALADIN-LAEF and AROME-EPS on 29 April 2014. a) Temporal evolution of the mean areal precipitation compared with INCA, and b) temporal evolution of the number of precipitation objects. Dashed and dotted lines represent the ensemble mean and grey shades the ensemble spread. c) Temporal evolution of S (structure), A (amplitude) and L (location) scores of the ensemble means of ALADIN-LAEF (black) and AROME-EPS (grey).