Atmospheric
Chemistry
and Physics
Discussions

1    **An Evaluation of the Efficacy of Very High Resolution Air-Quality**
2    **Modelling over the Athabasca Oil Sands Region, Alberta, Canada**

3    **Matthew Russell[1], Amir Hakami[1], Paul A. Makar[2], Ayodeji Akingunola[2], Junhua Zhang[2], Michael D. Moran[2], and**
4    **Qiong Zheng[2]**

5    [1]Department of Civil and Environmental Engineering, Carleton University, Ottawa, Canada

6    [2]Air Quality Research Division, Environment and Climate Change Canada, Toronto, Canada

7

8    ## Abstract

9    We examine the potential benefits of very high resolution for air-quality forecast simulations using a nested

10   system of the Global Environmental Multiscale – Modelling Air-quality and Chemistry chemical transport model.

11   We focus on simulations at 1km and 2.5km grid-cell spacing for the same time period and domain (the industrial

12   emissions region of the Athabasca Oil Sands).  Standard grid-cell to observation station pair analyses show no

13   benefit to the higher resolution simulation (and a degradation of performance for most metrics using this

14   standard form of evaluation).  However, when the evaluation methodology is modified, to include a search over

15   equivalent representative regions surrounding the observation locations for the closest fit to the observations, the

16   model simulation with the smaller grid cell size had the better performance.  While other sources of model error

17   thus dominate net performance at these two resolutions, obscuring the potential benefits of higher resolution

18   modelling for forecasting purposes, the higher resolution simulation shows promise in terms of better aiding

19   localized chemical analysis of pollutant plumes, through better representation of plume maxima.

20   ## 1    Introduction

21   Numerical modeling of the atmosphere in an Eulerian framework relies on discretization of the computational

22   domain into a numerical grid. The horizontal grid cell size of atmospheric simulations can range in from hundreds

23   of kilometers, to the metre-scale of Large Eddy Simulation models. For the purposes of this study, Very High

24   Resolution (VHR) modelling refers to chemical transport models (CTMs) employing a horizontal grid cell spacing of

25   1km or less.  It is in this regime that the photochemical processes may be forecasted with resolved microphysics

26   (e.g. Milbrandt and Yau, 2005(a,b)), and detailed particle and gas-phase chemistry, using currently available

27   computer technology. VHR modelling is very computationally expensive, and also introduces its own set of

28   challenges, such as the availability of surface boundary condition fields as the model grid cell size decreases.

29   Moreover, it is not currently clear whether decreases in model grid cell size leads to more accurate results when

30   compared to observations. The motivation behind VHR modelling in CTMs is to reduce the impact of diluting

31   chemical concentrations - especially from averaging emission plumes into large grid cells – in order to better

32   capture inhomogeneities in emission profiles, to better simulate local transport processes associated with terrain

33   that would otherwise be smoothed by the use of a coarse grid, and to reduce truncation errors and hence achieve

34    better numerical accuracy (Jacobson, 1999).

35    We note here that while the terms "grid cell size" and "resolution" tend to be used interchangeably in the

36    literature, this is not true in a precise mathematical sense; more formally, the ability to resolve features of size

37    $2\Delta x$ requires a grid cell spacing of size $\Delta x$, and the highest spatial frequency which can be reconstructed from a

38    discrete sampling of the latter grid cell spacing will be $\frac{1}{2\Delta x}$, the Nyquist wavenumber of the grid cell size

39    discretization. Furthermore, atmospheric models may make use of energy dissipation techniques that broaden

40    the size of resolvable wavelengths to $3\Delta x$ to $4\Delta x$ (Grasso, 2000; Pielke, 2001). Model resolution is thus a function

41    of, but not equivalent to, grid cell size. Here, we define "resolution" as the ability of a model to clearly distinguish

42    components of a predicted atmospheric variable, as a *function* of grid cell size.

43    The issue of a model to distinguish these features is also compounded by uncertainties in model inputs. For

44    example, in a large rural setting, a large model grid cell will represent an area containing many roads, whose

45    emissions will be averaged into one value per species per time. As the grid cell size decreases however, this

46    averaging effect will be reduced, giving each road's emissions more impact on the resulting concentrations in the

47    grid cell containing it. However, the smaller grid cell size will also result in steeper concentration gradients in the

48    model between adjacent grid cells, which can in turn result in numerical instabilities that contaminate predictions

49    (Salvador, et al., 1999). At the same time, a reduction in grid-cell size can be shown formally to reduce

50    inaccuracies in the discretization of the governing equations for atmospheric motion (Coiffier, 2011). Previous

51    efforts to address these issues through variable grid size or structure in air quality modeling have not received

52    sustained attention, and therefore most current air quality models use a uniform (albeit nested) grid cell size in

53    applications (Garcia-Menendez *et al.*, 2010; Kumar *et al.*, 1997).

54    As resolution increases further, the shape of local topographical features (*e.g.* buildings and street canyons)

55    become more important. Both the increased topographic complexity, and potential numerical instabilities can

56    lead to differences in meteorological forcing as resolution increases (Wolke, *et al.*, 2012; Gego, *et al.*, 2005)). The

57    contribution of meteorological uncertainties due to resolution become more significant, especially for secondary

58    pollutants such as ozone (Valari and Menut, 2008) or secondary Particulate Matter (PM). For example, Markakis *et*

59    *al.* (2015) in their analysis of 4 km CHIMERE simulations for the relatively flat terrain of Paris, France, suggested

60    that model meteorological grid cell size does not significantly impact forecast accuracy. That may not have been

61    the case, had their terrain been more complex. In contrast, Queen and Zhang (2008) observed considerable

62    meteorological sensitivity to the more complex terrain in their 4 km resolution Community Multiscale Air Quality

63    (CMAQ, EPA 1999) model simulations simulation over the Appalachian Mountains in the eastern United States.

64    A number of studies, employing various approaches, have tried to evaluate the benefits of higher resolution

65    simulations by quantifying sub-grid variability by employing larger model grid cell sizes (Vardoulakis *et al.*, 2003;

66    Ching *et al.*, 2006; Pepe *et al.*, 2016). These studies have often demonstrated that failure to account for higher

67    resolution features may result in mischaracterization of concentrations or health impacts (Isakov *et al*., 2007).

68    Population exposure studies using air pollution models may be affected by resolution in a more complex fashion,

69    given that both the predicted field (a pollutant with a known health impact) and the data to which the predicted

70    field is to be linked (the human population) both have resolution dependencies.

71    Terrain and meteorology are not the only factors that contribute to greater uncertainties as horizontal grid cell

72    size is reduced – for example, the ability of the model to locally resolve emission fluxes may also become a factor.

73    This may result in improved or deteriorated model performance as the size of the grid cells decrease. Gridded

74    model emissions may have an intrinsic resolution dependence in the underlying spatial disaggregation fields, and

75    this can contribute to uncertainties and errors in emissions as grid cell size is decreased. For instance, Valari and

76    Menut (2008) found that the discrepancy between their modelled and observed concentrations grew, rather than

77    shrank, in response to decreases in grid cell size from 48km to 6 km, and they associated these results with

78    changes in the resulting local emission fluxes. They showed that in their model setup, with regard to ozone, a grid

79    cell size was reached ($12 \times 12$ km$^2$) where errors in inputs (errors in the emission inventory, wind direction, *etc.*)

80    outweighed the importance of other sources of model error such as grid cell size. The authors however noted that

81    Paris' ozone photochemistry very often resides on the transition between a $NO_x^-$ sensitive and a VOC–sensitive

82    regime (Sillman *et al.*, 2003). These are chemical conditions which can alternatively produce or titrate ozone, and

83    hence has a degree of sensitivity to precursor emissions, and therefore, also, to any errors in those emissions.

84    Conversely, in a 3-level nested 9- to 3- to 1- km MM5–CMAQ simulation over Osaka, Japan, Shrestha *et al*., (2009)

85    found that ozone comparisons to observations improved as the grid resolution increased. This was also the case

86    for a 36- to 12- 4-km nested MM5–CMAQ simulation over Houston, USA (Ching *et al*., 2006), where the ozone

87    forecast improvement associated with higher resolution was attributed to the ability of the finer grid cell size

88    model nests to adequately resolve high concentrations of freshly emitted NOx and hence allow for more local

89    ozone titration. The latter process might not take effect until the grid cell size is sufficiently fine to resolve the $NO_x$

90    source patterns (*i.e.,* a level where traffic and industrial sources can be identified.) This titration was not seen until

91    they decreased their grid cell sizes to 2 km and smaller. Stroud *et al.* (2011) noted a similar grid cell size

92    dependent chemical impact on model performance, where secondary organic aerosol formation maxima were

93    better simulated with a 2.5km grid cell size model than a 10km grid cell size model. In general, the impact of

94    resolution on model performance appears to depend on a number of factors, such as the terrain, spatial

95    distribution of sources, pollutant of concern, season, *etc*. (Arunachalam *et al.,* 2006; Queen and Zhang, 2008; Dore

96    *et al.*, 2012).

97    Whether or not simulated quantities improve with reference to observed quantities as applications approach VHR

98    grid cell sizes, the resulting *distribution* of the quantities tends to be more physically realistic (Dore *et al*., 2012;

99    Salvador *et al*., 1999; Valari and Menut, 2008).

Atmospheric
Chemistry
and Physics
Discussions

100

101    The benefits for model performance with increased spatial resolution are unclear, based on the above literature.

102    However, most papers converge towards the following qualitative conclusions:

103        1. The impact of terrain topology on meteorological forcing as grid cell size decreases can dwarf the impact of

104           a more accurate spatial apportionment of the corresponding emissions.

105        2. Decreases in grid cell size result in a more realistic spatial distribution of chemical species, whether or not

106           model performance is improved.

107        3. Uncertainties of spatial and temporal emissions allocation have an increasing influence on overall model

108           uncertainty as model grid cell size decreases.

109    The 1980's saw several studies in which the potential impacts of wind direction errors on dispersion model

110    performance were examined.  Fox (1981) noted that pairing of model output at observation station locations could

111    be done as a function of both time and space, as a function of time (combining the data across all stations), as a

112    function of space (combining all times, at each station location), or without any pairing (observations and data are

113    compared as cumulative frequency distributions).  The accuracy of regulatory dispersion models in the early 1980's

114    was such that Fox (1984) concluded that model and observation values paired in time and space exhibited "little to

115    no correlation" and discussed potential errors associated with transport.  Poor correlations were also noted by

116    Hanha (1988), reporting on the first generation of reactive-transport models, stated "wind direction errors are the

117    major cause of the poor agreement in hourly predictions of concentrations at short distances downwind of point

118    sources," as well as describing metrics for air-quality model evaluation.  Hanha (1988) also noted that model

119    predictions could be offset in space and time relative to observations, leading to poor performance statistics,

120    despite a greater degree of similarity of behavior if the offsets are taken into account.  Errors in wind-field

121    modelling were described as the main source of error in simulations of plumes by Carhart $et$ $al$ (1989), again

122    showing how better agreement resulted when model and observations were unpaired in time and/or space, and

123    noted that other metrics such as maximum plume width might better represent model performance.  Lee (1987)

124    found that small perturbations in space and time could result in poor correlations, despite similar histogram

125    distributions of both model and observations.

126    More recently, Kang $et$ $al.,$ (2007) examined the concept of using the area of the limiting resolution of the model (2

127    to 3Δx, where Δx is the horizontal grid cell size) to weight or spatially average model evaluation metrics for a single

128    grid-cell size, noting how the model's rated ability to capture high concentration events ("hits") was increased when

129    the limiting resolution of the model was incorporated into the performance metrics.  However, the use of averaging

130    may mask the potential for a model with a small grid cell size to contain both the desired plume magnitude, as well
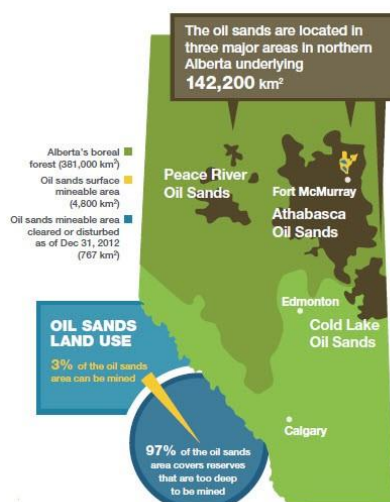
131    as much lower concentrations, within the same larger representative area, masking the potential impact of the

132    reduction in grid cell size.

133    We expand on this concept to evaluate the impact of model grid cell size in the context of an equivalent area about

134    a given observation location.  We examine area-weighted metrics in the form of averages over roughly equivalent

135    areas for different model grid cell sizes, and also use the *a priori* knowledge of the observations to determine

136    whether the closest match to observations may be found within an equivalent area. We show that the latter metric

137    demonstrates a positive impact of model grid cell size on simulation results, while more simple paired comparisons,

138    and averages over similar areas, mask these benefits.

139    We examine the impact of grid cell size on model performance in a region of intense petrochemical extraction and

140    upgrading, the Athabasca Oil Sands Region (AOSR). The AOSR refers to the northernmost of three large bitumen

141    deposits located the northern part of the province of Alberta in Canada; the Athabasca, Peace River, and Cold Lake

142    areas. Together these areas cover 142,200 $km^2$ in total, and constitute the third largest oil reserves in the world

143    (Government of Alberta, 2016), as shown in Figure 1.  The oil sands sector is the second largest source of $SO_2$ and

144    the third largest source of industrial $NO_x$ in the province of Alberta. This sector is also a significant source of

145    industrial PM, CO, and Volatile Organic Compound (VOC) emissions (Zhang *et al*., 2018), from a variety of source

146    types and industrial processes (*e.g.* open pit mine tailings ponds, large diesel fleets, bitumen upgrading facilities).

147    As is described below, very high resolution emissions data are available for these sources, and emissions take place

148    in a region with significant topography, hence the region provides a good test case for the relative impact of grid

149    cell size on air-quality model prediction results.

150    We describe next our model, the simulation domains and forecasting setup, the emissions data, our evaluation

151    methodology, and the results of our analysis.

152

153      Figure 1.  Map showing the Oil Sands regions (Government of Alberta, 2016).

154

Atmospheric
Chemistry
and Physics
Discussions

## 2 Methodology

### 2.1 GEM-MACH

The air-quality model used in this work is Environment and Climate Change Canada's (ECCC) Global Environmental Multiscale – Modelling Air-quality and Chemistry (GEM-MACH) model, which has been in use as Canada's operational air-quality forecast model since 2009 (Moran *et al.*, 2010). GEM-MACH is an on-line model, that is, both meteorological and chemistry processes are handled within a single model. The chemical processes reside within the physics module of the Global Environmental Multiscale meteorological forecast model (Côté, *et al.*, 1998(a,b)), originate with Environment Canada's earlier off-line model (A Unified Regional Air-quality Modelling System; AURAMS, Gong *et al.*, 2006), and include process representation for particle microphysics (Gong *et al.*, 2003(a,b)), inorganic heterogeneous chemistry (Makar *et al.*, 2003), aqueous phase chemistry, in-cloud and below-cloud scavenging (Gong *et al.*, 2006), and secondary organic aerosol formation (Stroud *et al*, 2011). GEM-MACH employs a sectional approach to represent the size distribution of atmospheric particles, with 12-bin (Makar *et al.,* 2015(a,b); Gong *et al.,* 2015) or 2-bin configurations (Moran *et al.*, 2010). The latter configuration is designed for maximum computational efficiency, with re-binning to the 12-bin distribution for key particle microphysics processes, in order to improve accuracy. Here, the 2-bin version of the model has been used, the main focus of the work being the impact of horizontal grid cell size on model results. Eight aerosol chemical components are resolved in GEM-MACH (sulphate, nitrate, ammonium, elemental carbon, primary organic aerosol, secondary organic aerosol, sea-salt and crustal material). In the present study, we make use of GEM-MACH v.1.5.1, described in more detail in Makar *et al*., 2015(a,b), employing 80 levels in a hybrid vertical coordinate system extending up to 0.1hPa (~30km).

### 2.2 Model Setup

#### 2.2.1 Grid Nesting

Four levels of nesting have been employed in our simulations, shown in Figure 2(a). This version of GEM-MACH operates on a rotated latitude-longitude coordinate system wherein the position of the coordinate system poles may be set by the user, allowing rotations of the grid with decreasing grid cell size during nesting. The outermost nested grid corresponds to the westernmost $^2/_3$ of the operational GEM-MACH forecasting domain, with a 10km grid cell size. Within that is nested a 10km grid cell size western Canada domain (yellow region, Figure 2(a)) which has been rotated to match the horizontal orientation of the Rocky Mountains, and which makes use of a similar double-moment microphysics scheme (Milbrandt and Yau, 2005 (a,b)) as the two innermost domains – this intermediate nested grid was constructed in order to allow hydrometeors to be passed from the western Canada 10km domain to the two innermost domains with a minimum of spin-up time required for the inner domain's meteorology. The third nested grid inwards (green region, Figure 2(a)) is the 2.5km grid cell size domain, which covers most of the Canadian provinces of Alberta and Saskatchewan. This grid will hereafter be referred to as the OS2.5km domain. The fourth and final nested grid (blue square, Figure 2(a)) is a 1km grid cell size domain, roughly centered over and covering the immediate environs of the Athabasca Oil Sands, and is referred to hereafter as the

Atmospheric
Chemistry
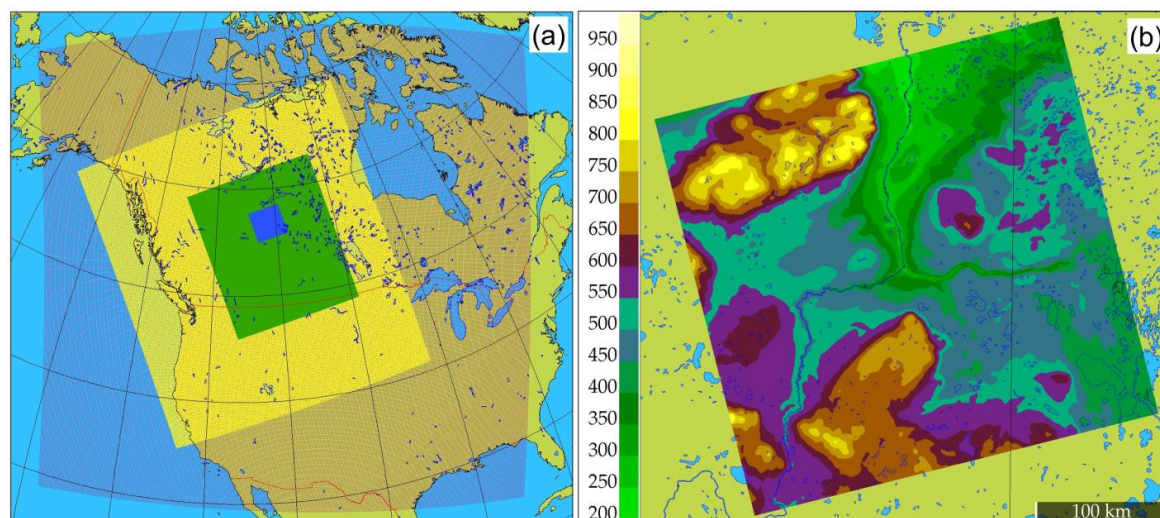and Physics
Discussions

191   OS1km model.  This last nest also shows the region within which 22 instrumented aircraft flights were conducted

192   during August and September of 2013, providing a unique measurement dataset for our evaluation of the

193   OS2.5km and OS1km model output for the same time period. Table 1 provides details on the horizontal

194   dimensions of each of these nested domains, and the duration of the simulations on each grid.  All four model

195   nests make use of the same vertical coordinate and levels.  Figure 2(b) shows the topography of the 1km domain

196   in detail; the region to be modelled is situated in a broad river valley, with a local vertical relief of 750 m.

197   Significant wind shears and frequent inversions are observed in the region, and part of our interest in 1km grid cell

198   size simulations is to determine the extent to which these local features may influence model prediction accuracy.

199   2.2.2 Simulation Cycling Strategy

200    The forecasts run in a repeating cycle from new meteorological analyses on every 36 hours, and hence are

201   constrained by observations to present chaotic drift of the forecast over an extended simulation. The outermost

202   10km domain carries out a 36 hour forecast, of which the first 6 hours is discarded as spin-up; the final 30 hours is

203   used as initial and boundary conditions for the rotated 10 km grid cell size domain (the OS10km domain).   As

204   noted above, the OS10km domain makes use of a microphysics package matching that of the subsequent higher

205   resolution simulations for better matching of cloud fields at those resolutions. An OS10km simulation of 30 hours

206   is then carried out, with the first 6 hours being discarded as spin-up, and the latter 24 hours forming the initial and

207   boundary conditions for the 2.5 km grid cell size OS2.5km simulation. The OS2.5km simulation is of 24 hours

208   duration.  The OS1km simulation covers the same 24 hours (and hence both 2.5km and 1km simulations start from

209   the same OS10km initial conditions at for every 24 hour forecast), with the 2.5km simulation providing boundary

210   conditions thereafter to the OS1km model. Continuity between 24 hour forecasts is thus maintained at the level of

211   the outermost nest.  The outermost domain is cycled every 12 hours starting at 0UT and 12UT; however, we have

212   selected the set of contiguous OS2.5km and OS1km 24 hour simulations starting from the 12UT continental

213   domain for our comparison.

214   Meteorological boundary conditions for lowest resolution GEM-MACH simulations are taken from operational

215   GEM forecasts, in turn driven by data assimilation analyses performed at the Canadian Meteorological Centre.

Atmospheric
Chemistry
and Physics
Discussions

216



217     Figure 2. (a) The four nested domains of the GEM-MACH simulations.  From outermost to innermost domains,

218     these are CONT10km (outermost, red dots), OS10km (yellow), OS2.5km (green), and OS1km (blue).  The model

219     simulations from the two innermost domains are the focus of the present study. (b) Topography in the OS1km

220     domain centred on Fort McMurray, Alberta (m agl).  The coloured area corresponds to the central blue domain in

221     (a).

222     Table 1.  Nested Domain Specifications

| Parameter | CONT10km | OS10km | OS2.5km | OS1km |
|---|---|---|---|---|
| Grid Size | 520x520 | 318x280 | 643x544 | 318x324 |
| Time step size (s) | 300 | 300 | 60 | 20 |
| Hours simulated | 36 | 30 | 24* | 24* |

223     *Note that both OS2.5km and OS1km output frequency was hourly.

224     2.3  Model Emissions

225     All emissions data used in this work are described in Zhang *et al*. (2018).  These emissions data include (a) direct

226     observations of stack-specific hourly emissions measured by Continuous Emission Monitoring Systems (CEMS), (b)

227     regional emissions inventory data from the Cumulative Environmental Management Association (CEMA) - which

228     had the most detailed  stack and process level emission data for the AOSR facilities, including emissions from mine

229     faces, tailings ponds, and the off-road mining fleet), (c) the 2010 Canadian Air Pollutant Emissions Inventory (APEI)

230     - which is the most comprehensive national emissions inventory, and which has the largest spatial coverage for

Atmospheric
Chemistry
and Physics
Discussions

231 area sources for areas outside the AOSR, and (d) the 2013 National Pollutant Release Inventory (NPRI) (a subset of

232 the APEI) that is based on emissions reports from large industrial facilities.

233 These emissions data sets primarily describe emissions of pollutants known as criteria-air-contaminants ($NO_x$,

234 VOCs, $SO_2$, $NH_3$, CO, $PM_{2.5}$, and $PM_{10}$) for *major-point sources* (*i.e.*, large emission stacks) and *area sources*. Area

235 emissions sources typically consist of multiple small mobile sources spread over a large area (*e.g.,* off-road

236 vehicles), large flux sources such as mine tailings settling ponds or mine faces, and/or large numbers of small

237 stacks for which no stack characteristic data (volume flow rates, temperatures of emissions, stack diameters),

238 needed to estimate plume-rise heights, are available.

239 Major-point sources are represented by a single geographical (latitude, longitude) pair of coordinates, and are

240 assigned to the grid cell in which the point is located. These sources are likely to be the most impacted by model

241 horizontal grid cell size, as even a large major-point source plume, which in reality may only occupy an emissions

242 area on the order of 100 $m^2$, is represented by a flux spread over an entire grid cell. A plume from a major point

243 source within a 2.5km grid cell will thus be immediately diluted to a size of 6.25$km^2$ upon emission, whereas the

244 same source with a 1km grid cell will have a cross-sectional horizontal extent of 1$km^2$. At the same time, higher

245 resolution may require a much more accurate representation of model winds close to the sources to maintain

246 accuracy in evaluation metrics dependant on plume position such as correlation – a wider plume being more likely

247 to at least partially intersect a monitoring station location than a narrower plume.

248 Area sources that are large compared to both model grid cell sizes (2.5km and 1km) can be expected to be

249 approximated by model grid cells of both resolutions, and are thus expected to be less impacted by model

250 resolution than emissions from point sources. However, smaller area sources (*i.e.* areas intermediate between

251 2.5km and 1km to the side) may be better resolved, and hence have less dilution and higher downwind

252 concentrations, when higher spatial resolution is employed.

253 In the AOSR, approximately 95% of the $SO_2$ emissions originate in major-point sources, while $NO_2$ is

254 approportioned ~40% to major-point sources and ~60% to area sources (Zhang *et al.*, 2018). Consequently our *a*

255 *priori* expectation is that the impact of the resolution change will be strongest for species like $SO_2$, less strong for

256 species like $NO_2$ that are emitted in part by point sources, but may also be apparent for other species and

257 secondary products, such as $O_3$.

258    2.4   Model Evaluation Methodology and Metrics

259 Comparisons between air-quality models and observations usually take the approach of comparing observation

260 and model-generated values paired in time and space, from the observation location and corresponding model

261 grid-cell respectively. We refer to this approach hereafter as our "standard" evaluation, for both 2.5km and 1km

262 simulations. However, we note additional factors aside from grid-cell size may influence the outcome of air-

Atmospheric
Chemistry
and Physics
Discussions

263    quality model evaluations.    For example, the relative skill of the meteorological component of the air-quality

264    model will depend in part on the density of meteorological observation data, incorporated into the model via data

265    assimilation, for the construction of the model's initial meteorological state.  This in turn will influence the local

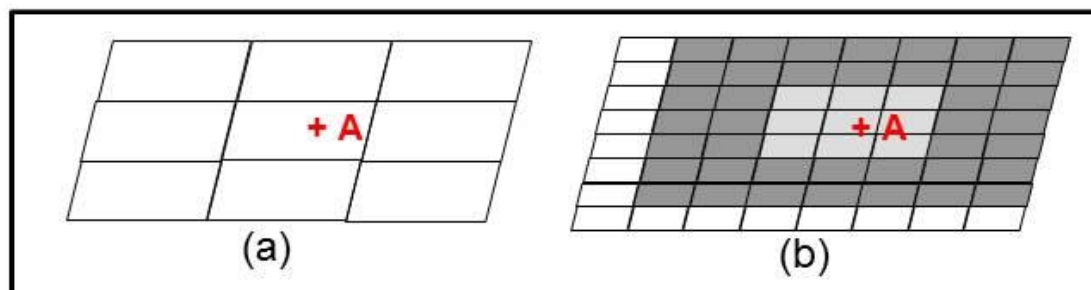266    skill of the model's predicted wind directions and hence the skill of its plume transport.  The simulations carried

267    out here focus on the Fort McMurray area, where the nearest available upper air meteorological sounding site is

268    located at the ECCC Stony Plain station, located approximately 500km south-west of the study area.    The

269    advantage of higher resolution simulations (*e.g.*, reduced numerical error associated with the discretization of

270    transport operators, and better treatment of local topographic influences) may thus be offset by errors in the

271    predicted *large scale* flow.

272    While meteorological model synoptic-scale forecast errors may manifest themselves locally as errors in the

273    direction of winds driving local plume transport, other advantages may result from the use of higher resolution

274    air-quality models.   Since lower resolution models *de facto* instantaneously redistribute plumes emitted from

275    large stack sources over a larger area, such artificial diffusion will reduce the model's ability to accurately simulate

276    concentration maxima, and the resulting chemistry, within simulated model plumes.   However, the spatial extent

277    of a plume in a model employing a large horizontal grid cell size may be such that its existence may be captured at

278    discrete observing sites.    In contrast, forecast plumes in models with smaller horizontal grid cell sizes may

279    correctly capture plume magnitude and chemical behaviour, but may be more subject to errors in the larger scale

280    wind direction.   To illustrate this point, Figure 3 shows a conceptual diagram of an actual plume, a large grid cell

281    size model plume, and a small grid cell size model plume, where the latter two simulated plumes are both subject

282    to the same synoptic-scale error in wind forecast direction (indicated by large red arrows; the smaller red arrow in

283    Figure 3(c) indicates the impact of local forcing predicted for the second model).  Observation station "+A" is

284    located downwind, and records the presence of the actual plume (Figure 3(a)).  The coarse grid cell size simulated

285    plume (Figure 3(b)), despite the error in the forecast wind direction, captures part of the observed plume in the

286    resulting time series at the observation station location.   In contrast, the small grid cell size plume (Figure 3(c)),

287    despite resolving the plume shape (and plume-internal chemistry) to a greater degree than the coarse grid cell size

288    simulated plume, fails to record the presence of the plume at the observation location.   A simple paired

289    observation-model time series evaluation would thus suggest that the former model has superior performance to

290    the latter model in this example, despite the latter model having created a more "realistic" plume in terms of the

291    maximum concentration reached, albeit in the wrong location, due to synoptic-scale forecast wind direction error.

292    In this particular instance, the magnitude of the smaller grid cell size simulated plume is more realistic than that of

293    the coarse grid cell size plume, but this improvement will not be captured in a standard evaluation analysis.  Shifts

294    in plume location across individual grid cells away from the location of an *in-situ* observation are more likely grid

295    cell size decreases.  In this example, a standard analysis would impose a more stringent expectation on the smaller

296    grid cell size simulation to correctly identify plume locations.

Atmospheric
Chemistry
and Physics
Discussions

297



298 Figure 3. Schematic comparison of surface concentration contours and model grid cell values of a transported pollutant

299 plume from a large stack (termed a "point" source). Wind direction shown by red arrows. Monitoring station location

300 marked by "+A". (a) Actual plume. (b) Coarse grid cell size air-quality model prediction. (c) Fine grid cell size air-quality model

301 prediction. Note the change in wind direction between observations (a) and simulations (b,c) associated with errors in the

302 forecast of the synoptic wind.

303 In order to attempt to evaluate the potential for higher resolution simulations to provide potential benefits that

304 are masked by synoptic scale forcing errors, in addition to the standard analysis, we perform additional analyses

305 that examine the model's ability to resolve plumes in the *vicinity* of the observation station. This is illustrated in

306 Figure 4.

307



308 Figure 4. Scale diagram of the same region in (a) 2.5km grid cell size simulation and (b) a 1km grid cell size simulation.

309 Region enclosed by light grey / dark grey shading in (b) represents the nearest nine / forty-nine 1km gridpoints surrounding

310 the observation location "A".

311 Figure 4(a) shows an observation station enclosing the nine nearest-neighbour model grid-cells for a 2.5km grid

312 cell size, while Figure 4(b) shows the corresponding 1 km grid cell size map, with the nine nearest-neighbour

313 model grid-cells shown in light grey, the forty-nine nearest grid cells shown in the region enclosed in dark grey.

314 Figure 4(a) encloses a region of 56.25 $km^2$ (7.5x7.5 km), while the light grey region in Figure 4(b) encloses 9$km^2$,

315 and the darker grey region encloses 49 $km^2$.

316 As noted above, in a formal mathematical sense, the smallest region resolvable by an Eulerian grid model is twice

317 the size of the model grid cell size (relating to the Nyquist frequency of the model); hence the smallest resolvable

318 feature spans two model grid cells in each direction. However, in a practical sense, a total of nine grid cells

319 centred on the observation station must be used to allow a boundary of two grid cells in any direction. Sampling

Atmospheric
Chemistry
and Physics
Discussions

320    any or all of the 9 grid cells in Figure 4(a) may thus be said to be representative of the model's ability to simulate

321    events occurring at discrete location "+A". The closest corresponding sampling region available to the 1 km model

322    (Figure 4(b)) is shown in dark grey. The light grey region of Figure 4(b) represents the closest 1 km grid cell size

323    region that corresponds to the single 2.5 km grid cell in which the observation station is located in Figure 4(a). We

324    attempt to ascertain model performance in these approximately equivalent regions around each observation

325    station, in the analysis that follows.

326    Our approach follows two steps:

327    (1) From the 2.5km simulation, in addition to the predicted model value at the grid-cell containing the

328        observation location, we determine the model grid-cell value in the nine grid-cells surrounding the

329        observation station location which has the closest value to that observed at the station. This represents the

330        model's "best estimate" of the value at the observation station location itself, to the model's ability to resolve

331        features at 2.5km grid cell size.

332    (2) From the 1km simulation, in addition to the model value at the grid-cell location, we select the closest value to

333        the observation value from: (a) the nearest nine grid-cells to the observation station location, and (b) the

334        nearest 49 grid-cells to the observation station location. The former represents the model's "best estimate"

335        of the value at the observation station location itself, while the latter represents the 1km model's best

336        estimate in the closest equivalent region to the limiting resolution of the 2.5km model.

337    Comparing the resulting statistical measures of each of these selected values with observations, in addition to the

338    standard analysis, thus evaluates the model's best attempt to resolve features for the specified grid cell size, and

339    allows cross-comparison of model performance within nearly equivalent areas. Cross-comparing the statistical

340    values for the different regions described above shows the model's ability to resolve features such as plumes from

341    the standpoint of the region represented at the different grid cell sizes. If synoptic-scale transport direction errors

342    creates situations similar to that depicted in Figure 3(a), a standard comparison of error would be expected to

343    show little benefit to higher resolution. However, the "best model estimate" comparisons would capture the

344    ability of the higher resolution model to more accurately simulate the magnitude of the plume, if not its spatial

345    location. Each of these selection procedures will be employed in the surface concentration comparisons which

346    follow.

347    We evaluate our model simulations against observations made at surface monitoring networks in the vicinity of

348    the Athabasca oil sands, and aboard an instrumented aircraft, the National Research Council of Canada Convair.

349    For the surface monitoring data, hourly time series of model output were matched to station time series using the

350    different strategies described above. For the aircraft observations, we extract model values through temporal and

351    spatial interpolation to the aircraft's position during the flights and only perform the standard analysis, as well as

Atmospheric
Chemistry
and Physics
Discussions

352    examining the behaviour of the two simulations along cross-sections corresponding to the flight paths.

353    Our statistical metrics for evaluation are common to many other air-quality applications, and were computed

354    using the 'modstat' function from the OpenAir R package (Carslaw and Ropkins, 2012). The statistics calculated

355    here include: mean bias (MB; perfect score: zero), mean absolute gross error (MGE; perfect score: zero),

356    normalised mean bias (NMB; perfect score: zero), normalised mean gross error (NMGE: perfect score: zero), root

357    mean squared error (RMSE; perfect score: zero), correlation coefficient (r, perfect score: unity), coefficient of

358    efficiency (COE: a perfect score is unity, a zero/negative score means the model is equivalent/less predictive

359    than the mean of the observations), and the index of agreement (IoA; perfect agreement is unity, and -1

360    indicates no agreement or little variability).

361    ## 3    Simulation Comparisons and Evaluation

362

363    ### 3.1  Model-to-model comparisons and averages

364    We begin a comparison of 2.5km and 1km grid cell size for specific events, and for averages across the 1km

365    domain, in order to provide a qualitative comparison of the differences in simulations for the two simulations, and

366    then continue with the quantitative comparison. Figure 5 compares OS2.5km (left column) and OS1km (right

367    column) simulation results for a cross-section located 0.2km from a major $SO_2$ emissions source at 0, 12 and 24

368    hours into a given simulation day.

369



370 Figure 5. Comparison of simulated $SO_2$ plume mixing ratios (ppbv) located 0.2km from a major point source, for OS2.5km

371 simulations (left column) and OS1km simulations (right column), at 0 (a,b), 12 (c,d), and 24 (e,f) hours into a 24 hour

372 simulation.

373 The model results are identical at hour 0 due to both the OS2.5km and OS1km models being initialized from the

374 OS10km data at this time (small differences in Figure5(a,b) are due to slight mis-matches in the cross-section

375 locations). Subsequent cross-sections show the OS1km model is capable of resolving both higher absolute mixing

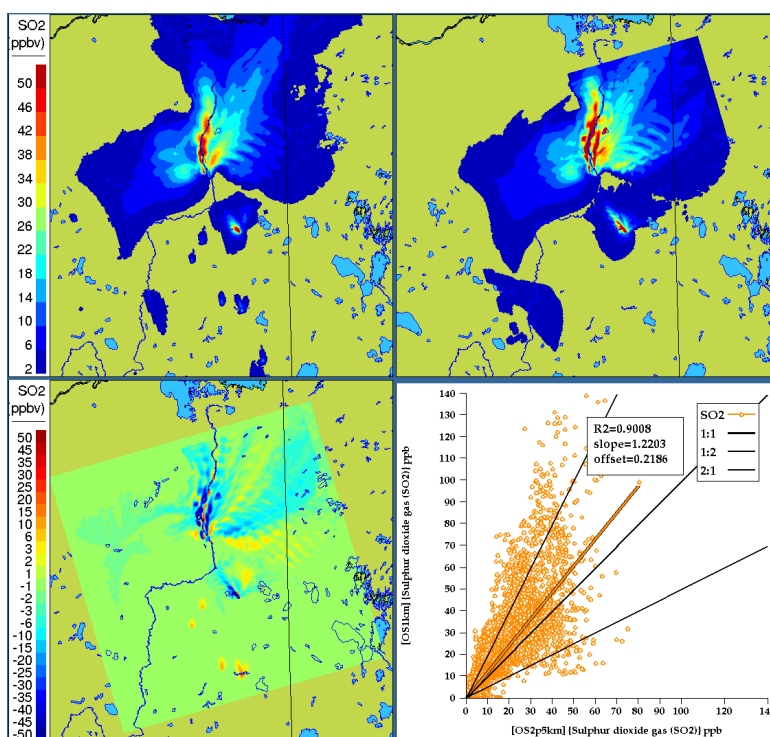376 ratio values, and sharper gradients, within 12 hours of simulation time (Figure 5 (c,d)). Multiple plumes are

377 resolved by 12 hours of simulation time in the 1km grid cell size simulation, along with markedly different plume

378 heights, plume structure, and a factor of two increase in the magnitude of plume mixing ratios relative to the

379 lower grid cell size simulation, and these differences persist into the 24[th] simulation hour (Figure 5(e,f)). Mixing

380 ratio differences of these magnitudes are to be expected given the increase in resolution, but Figure 5 shows that

381 other important aspects of the predicted plumes have changed. The plume heights are a function of predicted

382 local stability conditions in the grid-square containing the source, and the variation shown here represents a

383 substantial change in the predicted local stability for the origin sources of these plumes, resulting from the change

384     in model horizontal grid cell size.

385     Figure 6 compares the maximum surface $SO_2$ during the entire period for each simulation, as well as the difference

386     in maximum $SO_2$ between the simulations, along with a scatterplot of OS2.5km versus OS1km simulation results

387     (where, in the latter two panels, OS2.5km values were assigned to the corresponding OS1km grid-cell locations

388     using the nearest-neighbour approach).



Figure 6.  Comparison of total-simulation *maximum* surface $SO_2$ mixing ratios (ppbv) at (a) 2.5km and (b) 1km grid cell size (ppbv).   (c) Difference (2.5km – 1km).  (d) Scatterplot of 2.5km (x-axis) versus 1km (y-axis) total simulation average grid-cell surface $SO_2$ mixing ratios.

393     The maximum surface concentrations tend to show more elongated structures at the smaller grid cell size;

394     compare Figures 6(a,b), particularly for plumes in the western (left) half of the OS1km domain.  The difference

395     plot (Figure 6(c)) shows that local maximum concentration differences of up to -45 ppbv occur, due to changes in

396     the placement and maximum concentration of high concentration plumes.  The scatterplot of Figure 6(d) shows

397     that OS1km model has a demonstrated ability to achieve higher concentrations than the OS2.5km model, with a

398     slope of 1.22, and a noticeable clustering of values along the 1:2 line.  While these results are not unexpected

399     since approximately 95% of the $SO_2$ emissions in the domain originate in large stack, or point, sources, and hence

400     initial concentrations at source would be expected to 6.25x higher in the OS1km simulation, they also suggest that

401     a substantial improvement in the OS1km model's ability to capture $SO_2$ concentrations *should* be possible.  That is,

Atmospheric
Chemistry
and Physics
Discussions

402   the results of the two models are substantially different, and given the reduction in numerical error expected with

403   employing a smaller grid cell size, the latter might be expected to outperform a larger grid cell size model.

404   However, as we shall demonstrate in the next section, plume placement errors such as depicted in Figure 3 play a

405   substantial role in model performance as grid cell size decreases.

406   3.2 Quantitative comparisons
407
408   3.2.1 Surface observation comparison

409   The locations of the local network of 10 surface monitoring stations located near the sources of emissions in the

410   region (oil sands facilities) are shown in Figure 7.  As noted in section 2.4, we carry out several analyses:

411   (1) The standard evaluation (model values are extracted from the model grid-cells containing the observation

412       stations, at both grid cell sizes).

413   (2) Equal areas of representativeness, 1km and 2.5km grid cell sizes (the nearest nine OS1km grid cells are

414       compared to the OS2.5km single cell evaluation in two ways):

415       a.   Averaging of the OS1km results across the nine grid cells prior to evaluation (to determine whether

416            the mean value is better represented by the smaller grid cell size, similar to the approach taken in

417            Kang *et al.* (2007)).

418       b.   Selection of the *best* of the nine grid cells (closest to the observation value), to determine the extent

419            to which the OS1km model is capable of better representing the concentrations somewhere within

420            the corresponding OS2.5km model grid cell, if not at the OS1km cell closest to the observation

421            location.  Higher scores for the 1km grid cell size simulation in this case would indicate that while

422            errors in plume positioning (for example due to errors in the synoptic scale flow) negate some of the

423            advantages of the OS1km simulation, the plume may be better represented by the OS1km simulation

424            within the 2.5km grid cell's area.

425   (3) Equal areas of representativeness and equal regions of variability (nearest nine 2.5km cells are compared to

426       the nearest forty-nine 1km cells).  Here we make the assumption that the 2.5km grid cell size model's ability

427       to resolve features is limited to the surrounding three grid cells in each horizontal dimension, and make use of

428       the closest-in-size block of corresponding 1km cells (a $7 \times 7$ grid centered on the cell containing the

429       observation point.)  In both cases, the model value closest to the observations is chosen prior to evaluation.

430   While evaluations (2b) and (3) deliberately select the "best" value, they also provide a quantitative estimate of

431   the extent to which each model is capable of achieving the correct answer within roughly equal representative

432   areas centered on the observation station locations.  These comparisons are intended to evaluate (a) the

433   extent to which the 1km grid cell size is capable of improving simulation results despite, *e.g.*, the larger scale

Atmospheric
Chemistry
and Physics
Discussions

434    flow resulting in errors in the plume placement, and (b) whether the 1km grid cell size model is capable of

435    outperforming the 2.5km grid cell size model *over equivalent regions*.  In the last test, we place both models on

436    an equal footing with regards to the region being represented, as well with regards to allowing cell-to-cell

437    variability and the selection of a closest match to observations.

438    Our evaluation is presented as tables of statistical metrics.  The comparisons employing the nearest neighbour

439    approach are described with a "B#" superscript suffix, denoting that the "Best" sample within a square centred

440    on the observation point containing a total of # grid cells (*e.g.* the OS1km[B9] label denotes a comparison

441    between observed data and the simulation grid cell within a $3 \times 3$ grid-cell square centered about the

442    observation point).  Similarly, an A# superscript describes a comparison between the observations and the

443    Average of the # square of grid cells centered on the observation point.

444    Comparisons to surface concentrations were performed using publicly available data collected by the Wood

445    Buffalo Environmental Association (WBEA), which operates the air-quality monitoring network residing within

446    the OS1km domain. The monitoring station locations are shown in Figure 7.  The statistical performance of the

447    models, calculated using the procedure outlined above, are given in Tables 2 through 5, for $SO_2$, $NO_x$, $O_3$, and

448    $PM_{2.5}$, respectively.



449

450    Figure 7.  Illustration of the OS1km domain, with observation station locations. (a) Entire domain. (b) Close-up
451    view of station locations.  Monitoring stations are shown as purple outline squares in both images.  Light grey
452    regions in the background satellite image (b) are oil sands open-pit mining operations.

453    In the *standard* model grid cell to observation measurement comparison for $SO_2$, and $NO_x$ (first two columns,

454    Tables 2 and 3), the OS1km simulation had *worse* scores for all the metrics considered here.  For $O_3$, the OS1km

455    model had the better score for the correlation coefficient and root mean square error, and worse scores for all

18

Atmospheric
Chemistry
and Physics
Discussions

456    remaining model evaluation metrics. For $PM_{2.5}$, the OS1km model had higher performance for the correlation

457    coefficient and biases, while the OS2.5km model outperforms the OS1km model for all other metrics examined

458    here. Based on a standard analysis, the OS1km model thus performs poorly compared to the OS2.5km model; the

459    expected advantages associated with reduced numerical error in transport at smaller grid cell sizes are being offset

460    by other factors controlling the net model error.

461    When standard evaluation is compared to the *average* of the nearest nine 1km simulation grid cells surrounding

462    the observation point (first three columns of the tables), an intermediate result appears. For $SO_2$ (Table 2) the nine-

463    cell OS1km average has the best performance for correlation coefficient - indicating a better time distribution of

464    events may be achieved by a nine cell average at 1km grid cell size. The other metrics for the A9 simulation are

465    intermediate between the two standard evaluations for each simulation, indicating that some of the performance

466    loss resulting from the use of 1km grid cell size is reduced through averaging results to approximately the same size

467    regions as the OS2.5km grid cell size. The latter result holds for all metrics for $NO_x$ (including R, see Table 3). For

468    ozone (Table 4), averaging the nine nearest OS1km grid cells prior to measurement gives the best performance for

469    R and RMSE, and worse performance for the other metrics. For $PM_{2.5}$ (Table 5), all metrics for the OS1km nine grid-

470    cell average aside from the bias fall mid-way between the two standard methodology evaluations. Averaging the

471    smaller grid cell size model results thus shows a marginal improvement, depending on the species, but overall does

472    not compensate for the decrease in performance resulting from going to the smaller grid cell size.

473    We next ask the question, "Does a more accurate simulation value *exist* within the same region of the 1km model

474    as is encompassed by a 2.5km grid cell?" (fourth column of these Tables), by selecting the model value in the

475    nearest nine 1km grid cells with the closest match to observations and comparing to the corresponding single 2.5

476    grid cell. A dramatic improvement in the relative OS1km performance metric scores can be seen. For each of

477    Tables 2 through 5, this "best of nine" 1km comparison outperforms the previous 3 comparisons (columns 1

478    through 3), for all metrics. These improvements are sometimes dramatic (*e.g.* a doubling of correlation coefficient

479    along with a reduction in mean bias by a factor of three, a reduction of $NO_x$ mean bias values by a factor of 3, a shift

480    of coefficient of error from negative to positive values for $O_3$, and a reduction in the coefficient of error for $PM_{2.5}$ by

481    a factor of 2.5 compared to the nearest competing value from the previous evaluations. The coefficient of

482    efficiency for $SO_2$ and $O_3$ make the transition from negative to positive values when the "best-of-nine" methodology

483    is used, indicating that the model is able to better predict the observations than the observed mean, somewhere

484    within an equivalent area. This evaluation suggests that the OS1km model does *contain* a better result within the

485    same approximate region encompassed by a 2.5km grid cell. However, the location of that better result may be

486    subject to positioning error, such as described in Figure 3.

487    A valid argument could be made that the methodology employed in this fourth evaluation is subject to selection

488    bias, in that the selection of a *best* value in the case of the nearest nine 1km simulation places that model

Atmospheric
Chemistry
and Physics
Discussions

489    simulation at an advantage relative to the 2.5km model.  To address this last issue, the final two additional

490    methodologies for evaluation were employed, still maintaining the same approximate area of representativeness

491    for a grid cell, namely choosing the best value out of the nearest *nine* 2.5km grid cells (the limiting resolution of this

492    model simulation), and the best value out of the nearest *forty-nine* 1km grid cells (fifth and sixth columns of Tables

493    2 through 5, respectively).  That is, we attempt to place the two models on an equal basis with regards to selection

494    bias within a given region containing an observation station.

495    Two important results can be seen from this final evaluation.  First, as was the case for the "Best of 9" for the

496    OS1km simulation compared to the standard OS1km evaluation, the "Best of 9" for the OS2.5km simulation has a

497    considerably better performance than the standard OS2.5km evaluation (compare fifth and first columns, Tables 2

498    through 5). That is, the OS2.5km model may *also* be subject to location errors in transported species representation

499    which influence model performance.  However, when performance within the 56.25 $km^2$ area surrounding each

500    measurement point in the OS2.5km "Best of 9" evaluation is compared to the 49 $km^2$ area surrounding the

501    measurement points in the OS1km "Best of 49" simulation (*i.e.* compare columns five and six in Tables 2 through 5),

502    it can be seen that the OS1km model outperforms the OS2.5km model for all metrics for $O_3$, and $PM_{2.5}$, and all

503    metrics aside from bias for $SO_2$ and $NO_x$.   That is, despite the OS1km model having a slight disadvantage in the

504    relative size of the representative area containing the measurement station location, and both models being

505    allowed a similar selection strategy, the OS1km model is capable of generating values closer to the observations

506    than the OS2.5km model within an equivalent sub-region, across most of the metrics and chemical species

507    considered here.

508    This final result is strongly suggestive of the presence of issues such as illustrated in Figure 3.  These may include

509    errors in the larger scale synoptic wind flow, combined with the reduced size of plumes as grid cell size is reduced,

510    leading to more "misses" than "hits" for a given recorded event at a measurement station compared to the coarse

511    grid cell size model.  There may be multiple additional causes for such errors (examples include poor observation

512    density in the region for model initialization, underlying lower resolution boundary condition fields such as

513    topography not improving with the reduction in grid cell size, inaccuracies in land use fields used in meteorological

514    modelling due to rapid development, and errors in other aspects of the reaction transport modelling system aside

515    from horizontal resolution).  The expected advantages of the small grid cell size, such as better representation of

516    the concentrations of species within plumes and hence better representation of their reactive chemistry (c.f.

517    Lonsdale *et al.,* 2012), may be lost in a standard performance analysis due to these other issues.

518    Our analysis suggests that a practical limit in the benefits of increasing model accuracy may be reached when

519    resolution exceeds some threshold, as a result of other errors inherent in the modelling system.  However, the

520    analysis also suggests that if these non-resolution-related errors are corrected, the benefits of adopting a smaller

521    grid cell size may be substantial.  For example, meteorological data assimilation employing a dense monitoring

522 network for a specific area of interest would be expected to show a greater impact for smaller than larger grid cell

523 sizes, due to the greater ability of the former to take advantage of the observation density in correcting the initial

524 meteorological state. We note that recent work applying land use data assimilation (Carrera *et al.*, 2015) to

525 regional 2.5km grid cell size weather simulations (Milbrandt *et al.,* 2016) have suggested that such data assimilation

526 may indeed improve forecast skill at the very local scale.

527 Table 2. Surface $SO_2$ observations to model comparison for entire simulation period (ppbv)

| Evaluation Metric | OS2.5km | OS1km | OS1km$^{A9}$ | OS1km$^{B9}$ | OS2.5km$^{B9}$ | OS1km$^{B49}$ |
|---|---|---|---|---|---|---|
| IoA | 0.237 | 0.154 | 0.207 | 0.601 | 0.701 | 0.810 |
| r | 0.290 | 0.230 | 0.295 | 0.604 | 0.672 | 0.848 |
| NGME | 2.128 | 2.363 | 2.212 | 1.114 | 0.834 | 0.529 |
| GME | 2.918 | 3.240 | 3.034 | 1.528 | 1.143 | 0.725 |
| CoE | -0.525 | -0.693 | -0.585 | 0.202 | 0.403 | 0.621 |
| RMSE | 7.063 | 9.665 | 7.876 | 4.436 | 3.671 | 2.618 |
| NMB | 1.130 | 1.376 | 1.299 | 0.347 | -0.010 | 0.017 |
| MB | 1.550 | 1.887 | 1.781 | 0.475 | -0.013 | 0.024 |

528 • 5466 Samples used

529 Table 3. Surface $NO_x$ observations to model comparison for entire simulation period (ppbv)

| Evaluation Metric | OS2.5km | OS1km | OS1km$^{A9}$ | OS1km$^{B9}$ | OS2.5km$^{B9}$ | OS1km$^{B49}$ |
|---|---|---|---|---|---|---|
| IoA | 0.177 | 0.138 | 0.152 | 0.416 | 0.589 | 0.665 |
| r | 0.143 | 0.114 | 0.116 | 0.165 | 0.305 | 0.388 |
| NGME | 1.520 | 1.593 | 1.567 | 1.079 | 0.760 | 0.619 |
| GME | 12.898 | 13.518 | 13.296 | 9.156 | 6.447 | 5.255 |
| CoE | -0.646 | -0.725 | -0.697 | -0.168 | 0.177 | 0.329 |
| RMSE | 28.052 | 35.197 | 34.644 | 25.782 | 15.315 | 13.704 |

Atmospheric
Chemistry
and Physics
Discussions

| NMB | 0.493 | 0.570 | 0.542 | 0.174 | -0.027 | -0.063 |
|-----|-------|-------|-------|-------|--------|--------|
| MB  | 4.183 | 4.834 | 4.597 | 1.477 | -0.231 | -0.531 |

530
- 3257 Samples used

531

532     Table 4. Surface $O_3$ observations to model comparison for entire simulation period (ppbv)

| Evaluation Metric | OS2.5km | OS1km | OS1km[A9] | OS1km[B9] | OS2.5km[B9] | OS1km[B49] |
|-------------------|---------|-------|-----------|-----------|-------------|------------|
| IoA               | 0.414   | 0.405 | 0.404     | 0.527     | 0.637       | 0.690      |
| r                 | 0.496   | 0.506 | 0.515     | 0.606     | 0.688       | 0.738      |
| NGME              | 0.660   | 0.670 | 0.672     | 0.534     | 0.410       | 0.349      |
| GME               | 10.757  | 10.915| 10.949    | 8.692     | 6.673       | 5.687      |
| CoE               | -0.172  | -0.189| -0.193    | 0.053     | 0.273       | 0.380      |
| RMSE              | 16.040  | 15.859| 15.794    | 13.305    | 11.084      | 9.719      |
| NMB               | 0.527   | 0.559 | 0.579     | 0.463     | 0.337       | 0.304      |
| MB                | 8.579   | 9.104 | 9.431     | 7.536     | 5.488       | 4.945      |

533
- 2189 Samples used

534     Table 5. Surface $PM_{2.5}$ observations to model comparison for entire simulation period ($\mu g\ m^{-3}$)

| Evaluation Metric | OS2.5km | OS1km | OS1km[A9] | OS1km[B9] | OS2.5km[B9] | OS1km[B49] |
|-------------------|---------|-------|-----------|-----------|-------------|------------|
| IoA               | 0.280   | 0.262 | 0.267     | 0.412     | 0.508       | 0.572      |
| r                 | 0.201   | 0.216 | 0.214     | 0.314     | 0.376       | 0.466      |
| NGME              | 0.791   | 0.811 | 0.806     | 0.647     | 0.541       | 0.471      |
| GME               | 5.342   | 5.478 | 5.441     | 4.365     | 3.651       | 3.181      |
| CoE               | -0.439  | -0.476| -0.466    | -0.176    | 0.016       | 0.143      |
| RMSE              | 8.286   | 8.786 | 8.663     | 7.117     | 6.169       | 5.690      |

Atmospheric
Chemistry
and Physics
Discussions

| NMB | -0.268 | -0.257 | -0.257 | -0.289 | -0.299 | -0.287 |
| MB | -1.812 | -1.734 | -1.736 | -1.948 | -2.016 | -1.937 |

535             •     3377 Samples used

536 The surface observation data were also analyzed by time-of-day, with both observations and simulations split into

537 daytime (hours 9:00 to 18:00 local time) and nighttime (hour 19:00 to 8:00 local time) data pairs (Appendix, Tables

538 A1 through A8). Within each of these diurnally segregated time periods, the broad aspects of the comparison were

539 the same as for the "all data" Tables 2 to 5 above: the OS1km simulations tendied to have reduced performance in

540 a standard analysis, averaging improved but not completely ameliorated the performance of the OS1km simulation,

541 a methodology employing the best of nine OS1km grid cells had superior performance to the two standard

542 comparisons, and comparison of the "best of" methodologies for equal areas showed better performance for the

543 OS1km compared to the OS2.5km simulation. We also noted substantial differences in the day and night

544 performance of both models across the methodologies. For example, daytime $SO_2$ and $NO_x$ performance within a

545 given model and comparison methodology was usually better than nighttime performance for IOA,R, NGME, COE

546 and NMB, while worse for RMSE, while nighttime $O_3$ performance was better for IOA, r, NGME, and COE. Daytime

547 $PM_{2.5}$ performance was better than nighttime for IOA, r, COE, and NMB.

548 3.2.2 Comparisons to Aircraft Observations

549 Twenty-two aircraft observation flights were carried out during the study simulation period – we present

550 statistical comparisons using the standard approach only, here (model grid cell containing the observation point to

551 observation data at the aircraft location). Model values were linearly interpolated in time and space to the

552 aircraft observation locations and times (aircraft observations were on a 10s interval.) We begin with a composite

553 comparison across all observation times, in Table 6.

554 Table 6. Aircraft observation comparisons, $SO_2$ and $NO_2$ (ppbv)

|  | $SO_2$ (21787 samples) | | $NO_2$ (18310 samples) | |
| --- | --- | --- | --- | --- |
|  | OS2.5km | OS1km | OS2.5km | OS1km |
| IoA | 0.63 | 0.62 | 0.61 | 0.58 |
| r | 0.26 | 0.28 | 0.39 | 0.34 |
| NGME | 1.07 | 1.09 | 0.90 | 0.96 |
| GME | 3.98 | 4.06 | 1.56 | 1.68 |
| CoE | 0.27 | 0.25 | 0.23 | 0.17 |
| RMSE | 12.84 | 13.97 | 3.12 | 3.62 |
| NMB | -0.31 | -0.29 | -0.26 | -0.20 |
| MB | -1.17 | -1.07 | -0.45 | -0.34 |

555

556    The results are in general similar to the surface analysis, in that the OS1km simulation tended to have worse

557    performance than the OS2.5km simulation (exceptions being the biases for both $SO_2$ and $NO_2$, and the slightly

558    better OS1km correlation coefficient for $SO_2$).  One striking difference between the first two columns of Tables 2

559    and 3 and Table 14 are the magnitude of the differences between the simulations.  Aloft (Table 6), the differences

560    in performance metric magnitudes between OS2.5km and OS1km simulations are much smaller than at the

561    surface (Tables 3 and 4).  The biases are negative aloft, while positive at the surface, indicating that both models

562    may be lofting plumes to insufficient distances; one of the possible (non-horizontal grid cell size dependent)

563    causes of model error may be in the extent of vertical transport (this possibility is examined in more detail in

564    Akingunola *et al.,* 2018, and Gordon *et al.,* 2018).  An example of this behaviour is shown in Figure 8; both plumes

565    fumigate to the surface, while the observed plume resides largely aloft.  The OS1km model captures the higher

566    concentrations to a better degree, but the impact of excessive fumigation more than offsets this improvement, as

567    is shown by the performance evaluation of Table 7, where both models have negative biases aloft.  In this

568    particular case, the tendency of the model to overestimate the extent of fumigation has a bigger impact on

569    performance than grid cell size.



570

571    Figure 8.   Comparison between OS2.5km (top row) and OS1km (bottom row) simulations for $SO_2$ relative to
572    aircraft observations (ppbv).  (a,c): Simulated surface concentrations of $SO_2$, with the flight track shown as a red
573    line.   (b,d) Simulated concentration profiles along the flight path as a function of time, with the successive

Atmospheric
Chemistry
and Physics
Discussions

574   intersections of the flight path with the plume as background colour contours.  Observed SO$_2$ aboard the aircraft
575   are shown between the two black lines, which show the elevation of the aircraft on successive passes around the
576   facility.  Dotted lines show the upper and lower vertical extent of the observed plume.  Note that for both model
577   simulations, the plume erroneously fumigates the surface.

578

Atmospheric
Chemistry
and Physics
Discussions

579     Table 7.  Standard performance evaluation of Flight 8 for SO$_2$ (ppbv)

|        | OS2.5km | OS1km |
|--------|---------|-------|
| IoA    | 0.69    | 0.68  |
| r      | 0.42    | 0.31  |
| NGME   | 1.04    | 1.09  |
| GME    | 4.02    | 4.25  |
| CoE    | 0.39    | 0.35  |
| RMSE   | 16.72   | 20.57 |
| NMB    | -0.42   | -0.34 |
| MB     | -1.63   | -1.32 |

580                                   1261 samples used.


581     Meanwhile other flights show a clear advantage of the OS1km model.  One example is given by the NO$_2$

582     performance evaluation of Table 8 and depicted in Figure 9, for Flight 17 (a similar flight plan carried out around

583     the same facility as Flight 8).  While the correlation coefficient degraded slightly in the OS1km resolution

584     simulation, all other performance measures were improved with the decrease in grid cell size.  Two time versus

585     height profile cross-sections for Flight 17 are shown in Figure 9.  In the upper two panels, the OS2.5km (Figure

586     9(a)) and OS1km (Figure 9(b)) simulations are compared for the portion of the overall flight track circling the given

587     facility.  This comparison clearly shows that the OS1km model does a better job of capturing the width of the high

588     concentration region of the plume – however, the location of the model plume lags the observations.  During this

589     portion of the flight alone, the OS2.5km model statistics, particularly the correlation coefficient, outperform the

590     OS1km model, due to this issue of plume location mismatching.  Figures 9(a,b) may be compared to Figure 3(a,b) –

591     the same situation is depicted in both Figures.  Figure 9(c,d) show the OS2.5km simulation (10(c)) and OS1km

592     simulation results in another portion of the flight – here the OS1km performance for most statistics was better

593     than the OS2.5km model performance.  The OS1km model (Figure 9(d)) captures the existence of a lower

594     concentration layer aloft in the right-hand side of the cross-section, and the existence of low concentration

595     intervening layers, as well as the overall lower concentrations of SO$_2$, while the OS2.5km model does not resolve

596     these fine scale and lower concentration features.  We note here that IoA, CoE and the other error measures

597     capture the visual impression that the OS1km model outperforms the OS2.5km model for this flight, while the

598     correlation coefficient is highly dependent on the placement of the plume maximum in the upper two panels.


599     These and the snap-shot comparisons described in Section 3.1 show that the higher resolution model is having a

600     significant impact on predictions – however, other aspects of the overall model performance are preventing the

601     potential benefits of higher resolution from influencing the standard performance evaluation.


602


603

Atmospheric
Chemistry
and Physics
Discussions

604       Table 8.  Standard performance evaluation of Flight 17 for $NO_2$ (ppbv)

|  | OS2.5km | OS1km |
|---|---|---|
| IoA | 0.26 | 0.58 |
| r | 0.26 | 0.25 |
| NGME | 2.03 | 1.15 |
| GME | 0.52 | 0.29 |
| CoE | -0.48 | 0.16 |
| RMSE | 1.37 | 0.70 |
| NMB | 0.83 | -0.54 |
| MB | 0.21 | -0.14 |



605

606       Figure 9.  Flight 17 comparison for $NO_2$ (ppbv) for portions of the net flight track circling the CNRL facility for
607       OS2.5km (a) and OS1km (b) simulations, and for a later section of the same flight path for the OS2.5km (c) and
608       OS1km (d) simulations.

609

Atmospheric
Chemistry
and Physics
Discussions

## 4. Summary and Conclusions

Our work suggests the following:

Decreases to air-quality model horizontal grid cell size will not necessarily result in improvements to model performance in standard performance evaluations, in which the model values at the grid-cells encompassing measurement location stations are used in a pairwise comparison to observations. Other considerations, such as the accuracy of the larger scale wind direction and speed forecast, and the accuracy of the plume rise parameterization used within the model may play a greater role in the overall performance of the model, and reduce the benefits of the smaller grid cell size. In the context of a standard model performance evaluation, there may be fixed limits to the benefits of decreasing model grid cell size.

Despite this difficulty, our results also show that the use of smaller grid cell sizes have some potential benefits, in that these models do a better job of resolving specific air pollution features, like high concentration maxima within plumes. Both coarse and fine grid cell size plumes may be misplaced in both time and space, with the net result that the latter model has a worse performance in a standard comparison, but is nevertheless more likely to capture the correct in-plume concentrations, and hence the chemistry, of the actual plume, in the *neighbourhood* of the observation location. When the evaluation is broadened to find the closest fit to observations in the vicinity of the observation station, with models confined to a similar representative area around the observation station, these potential benefits of the smaller grid cell size become apparent.

These findings suggest that at the current state of development, VHR air-quality models are of benefit for the specific purpose of chemical process studies, in which the main aim of the work is to accurately simulate plume chemistry – and in which accurate forecasting of the *position* of the plume in time and space is a secondary concern. Our work also suggests that efforts to improve other aspects of the overall modelling framework which improve the large scale flow (for example, the use of data assimilation of local meteorology to improve wind direction predictions) may result in greater benefits as smaller grid cell sizes are employed.

Atmospheric
Chemistry
and Physics
Discussions

# 5   References

Akingunola, A., Makar, P.A., Zhang, J., Darlington, A., Li, S.-M., Gordon, M., Moran, M.D., Zheng, Q., A chemical transport model study of plume rise and particle size distribution for the Athabasca oil sands, *Atmos. Chem. Phys.*, 18, 8667-8688, 2018.

Arunachalam, S., Holland, A., Do, B. & Abraczinskas, M., A quantitative assessment of the influence of grid resolution on predictions of future-year air quality in North Carolina, USA. *Atm. Env.*, 40, 5010-5026, 2006.

Carhart, R.A., Policastro, A.J., Wastag, M., and Coke, L.,  Evaluation of eight short-term long-range transport models using field data, *Atm. Env.* 23, 85-105, 1989.

Carrera, M.L., Belair, S., Bilodeau, B.,  The Canadian Land Data Assimilation System (CALDAS):  Description and Synthetic Evaluation Study, *J. Hydromet.*, 16, 1293-1314, 2015.

Carslaw, D. C. and Ropkins, K., openair – an R package for air quality data analysis, *Environ. Modell. Softw.,* 27–28, 52–61, 2012.

Ching, J., Herwehe, J. and Swall, J., On joint deterministic grid modeling and sub-grid variability conceptual framework for model evaluation, *Atm. Env.*, 40, 4935-4945, 2006.

Coiffier, J.,  Fundamentals of Numerical Weather Prediction, Cambridge University Press, 363pp., 2011.

Côté, J., Gravel, S., Méthot, A., Patoine, A., Roch, M., Staniforth, A., The operational CMC--MRB global environmental multiscale (GEM) model. Part I: Design considerations and formulation, *Mon. Wea. Rev.*, 126, 1373-1395, 1998.

Côté, J., Desmarais, J.-G., Gravel, S., Méthot, A., Patoine, A., Roch, M., Staniforth, A., The operational CMC--MRB global environmental multiscale (GEM) model. Part II: Results. *Mon. Wea. Rev.*, 126, 1397-1418, 1998.

Dore, A. J., Kryza, M., Hall, J.R., Hallsworth, S., Keller, V.J.D., Vieno, M., and Sutton, M.A., The influence of model grid resolution on estimation of national scale nitrogen deposition and exceedance of critical loads. *Biogeosci.*, 9, 1597-1609, 2012.

EPA, 1999:  https://www.cmascenter.org/cmaq/science_documentation/ , last accessed September 2, 2018.

Fox, D.G., Judging air quality model performance - summary of the AMS Workshop on Dispersion Model Performance, Woods Hole, Mass., 8-11 September 1980, *Bull. Am. Met. Soc.*, 62, 599-609, 1981.

Fox, D.G., Uncertainty in air quality modelling – a summary of the AMS Workshop on Quantifying and Communicating Model Uncertainty, Woods Hole, Mass., September 1982, *Bull. Am. Met. Soc.*, 65, 27-36, 1984.

Garcia-Menendez, F., Yano, A., Hu, Y. and Odman, M. T., An adaptive grid version of CMAQ for improving the resolution of plumes. *Atm. Poll. Res.*, 1, 239-249, 2010.

Gego, E., Hogrefe, C., Kallos, G., Voudouri, A., Irwin, J.S., Rao, S.T., Examination of model predictions at different

Atmospheric
Chemistry
and Physics
Discussions

horizontal grid resolutions. *Env. Fluid Mech.*, 5, 63-85, 2005.

Gong, W., Dastoor, A.P., Bouchet, V.S., Gong, S.L., Makar, P.A., Moran, M.D., Pabla, B., Menard, S., Crevier, L-P., Cousineau, S., Venkatesh, S., Cloud processing of gases and aerosols in a regional air quality model (AURAMS), *Atm. Res.* 82, 248-275, 2006.

Gong, W., Makar, P.A., Zhang, J., Milbrandt, M., Gravel, S., Hayden, K.L., MacDonald, A.M., Leaitch, W.R., Modelling aerosol--cloud--meteorology interaction: A case study with a fully coupled air quality model (GEM-MACH). *Atm. Env.*, 115, 695-715, 2015.

Gong, S.L., Barrie, L.A., Lazare, M., Canadian Aerosol Module (CAM): a size-segregated simulation of atmospheric aerosol processes for climate and air quality models: 2. Global sea-salt aerosol and its budgets. *J. Geophys. Res*. 107, 4779. http://dx.doi.org/10.1029/2001JD002004, 2003a.

Gong, S. L., Barrie, L.A., Blanchet, J.-P., von Salzen, K., Lohmann, U., Lesins, G., Spacek, L., Zhang, L.M., Girard, E., Lin, H., Leaitch, R., Leighton, H., Chylek, P., Huang, P., Canadian Aerosol Module: A size-segregated simulation of atmospheric aerosol processes for climate and air quality models 1. Module development. *J. Geophys. Res.*, 108, D1, 4007, doi:10.1029/2001JD002002, 2003b.

Gordon, M., Makar, P.A., Staebler, R., Zhang, J., Akingunola, A., Gong, W., Li, S.-M., A comparison of plume rise algorithms to stack plume measurements in the Athabasca oil sands, *Atm. Chem. Phys. Disc.,* (https://www.atmos-chem-phys-discuss.net/acp-2017-1093/), 2018.

Government of Alberta, 2016: Alberta Energy: Oil Sands, http://www.energy.alberta.ca/oilsands/oilsands.asp, 2016, last accessed November 11, 2017.

Grasso, L.D., The differentiation between grid spacing and resolution and their application to numerical modelling, *Bull. Am. Met. Soc.*, 81, 579-580, 2000.

Hanha, S.R., Air quality model evaluation and uncertainty. *J. Air Poll. Cont. Assoc.*, 33, 406-412, 1988.

Isakov, V., Irwin, J. S., Ching, J., Using CMAQ for exposure modeling and characterizing the subgrid variability for exposure estimates. *J. App. Met. Cli*m., 46, 1354-1371, 2007.

Jacobson, M.Z., Fundamentals of Atmospheric Modelling, Cambridge U. Press, 656pp., 1999.

Kang, D., Mathur, R., Schere, K., Yu, S., Eder, B., New categorical metrics for air quality model evaluation, *J. App. Met. Clim.,* 46, 549-555, 2007.

Kumar, N., Russell, A.G., Segall, E., Steenkiste, P. Parallel and Distributed Application of an Urban-to-Regional Multiscale Model. *Comp. Chem. Eng.*, 21, 399-408, 1997.

Lee, I.Y., Numerical simulations of cross-Appalachian transport and diffusion. *Bound. Lay. Met.*, 39, 53-66, 1987.

Lonsdale, C.R., Stevens, R.G., Brock, C.A., Makar, P.A., Knipping, E.M., and Pierce J.R., The effect of coal-fired power-plant SO2 and NOx control technologies on aerosol nucleation in the source plumes, *Atm. Chem. Phys.*, 12, 11519-11531, 2012.

Atmospheric
Chemistry
and Physics
Discussions

Makar, P. A., Bouchet, V. S. & Nenes, A., Inorganic chemistry calculations using HETV--a vectorized solver for the SO42--NO3--NH4+ system based on the ISORROPIA algorithms. *Atm. Env.*, 37, 2279-2294, 2003.

Makar, P.A., Gong, W., Milbrandt, J., Hogrefe, C., Zhang, Y., Curci, G., Zabkar, R., Im, U., Balzarini, A., Baro, R., Bianconi, R., Cheung, P., Forkel, R., Gravel, S., Hirtl, H., Honzak, L., Hou, A., Jimenz-Guerrero, P., Langer, M., Moran, M.D., Pabla, B., Perez, J.L., Pirovano,G., San Jose,R., Tuccella, P., Werhahn, J., Zhang, J., Galmarini, S., Feedbacks between air pollution and weather, part 1: Effects on weather. *Atm. Env.,* 115, 442-469, 2015(a).

Makar, P.A., Gong, W., Hogrefe, C., Zhang, Y., Curci, G., Zabkar, R., Milbrandt, J., Im, U., Balzarini, A., Baro, R., Bianconi, R., Cheung, P., Forkel, R., Gravel, S., Hirtl, H., Honzak, L., Hou, A., Jimenz-Guerrero, P., Langer, M., Moran, M.D., Pabla, B., Perez, J.L., Pirovano,G., San Jose,R., Tuccella, P., Werhahn, J., Zhang, J., Galmarini, S., Feedbacks between air pollution and weather, part 2: Effects on chemistry. *Atm. Env.,* 115, 499-526, 2015(b).

Markakis, K., Valari, M., Perrussel, O., Sanchez, O., and Honore, C., Climate-forced air-quality modeling at the urban scale : sensitivity to model resolution, emissions and meteorology. *Atm. Chem. Phys.*, 15, 7703-7723, 2015.

Milbrandt, J. A. and Yau, M. K., A multimoment bulk microphysics parameterization, Part I: analysis of the role of the spectral shape parameter, *J. Atmos. Sci.*, 62, 3051–3064, 2005(a).

Milbrandt, J. A. and Yau, M. K., A multimoment bulk microphysics parameterization, Part II: a proposed three-moment closure and scheme, *J. Atmos. Sci.*, 62, 3065–3081, 2005(b).

Milbrandt, J.A., Belair, S., Faucher, M., Vallee, M., Carrera, M.L., and Glazer, A., The Pan-Canadian high resolution deterministic prediction system, Weather and Forecasting, 31, 1791-1816, 2016.

Moran, M. D. Ménard, S., Talbot, D., Huang, P., Makar, P. A., Gong, W., Landry, H., Gravel, S., Gong, S., Crevier, L.-P., Kallaur,A., Sassi, M., Particulate-matter forecasting with GEM-MACH15, a new Canadian air-quality forecast model. Air pollution modelling and its application XX. Springer, Dordrecht, pp. 289-292, 2010.

Pepe, N., Pirovano, G., Lonati, G., Balzarini, A., Toppetti, A., Riva, G.M., and Bedogni, M., Development and application of a high resolution hybrid modelling system for the evaluation of urban air quality. *Atm. Env.*, 141, 297-311, 2016.

Pielke, R.A. Sr., Further comments on "The differentiation between grid spacing and resolution and their application to numerical modeling", *Bull. Am. Met. Soc.*, 82, 699, 2001.

Queen, A. and Zhang, Y., Examining the sensitivity of MM5--CMAQ predictions to explicit microphysics schemes and horizontal grid resolutions, Part III—The impact of horizontal grid resolution. *Atm. En*v., 42, 3869-3881, 2008.

Salvador, R., Calbó, J. & Millán, M. M., Horizontal grid size selection and its influence on mesoscale model simulations. *J. App. Met.*, 38, 1311-1329, 1999.

Shrestha, K. L., Kondo, A., Akikazu, K. A. G. A.,  Inoue, Y., High-resolution modeling and evaluation of ozone air quality of Osaka using MM5-CMAQ system. *J. Env. Sci.*, 21, 782-789, 2009.

Sillman, S., Vautard, R., Menut, L. & Kley, D., O3-NO x-VOC sensitivity and NO x-VOC indicators in Paris: Results from models and Atmospheric Pollution Over the Paris Area (ESQUIF) measurements. *J. of Geophy. Res.*, 108, 8563, doi:10.1029/2002JD001561, 2003.

Stroud, C.A., P.A. Makar, M.D. Moran, W. Gong, S. Gong, J. Zhang, K. Hayden, C. Mihele, and J.R. Brook, Impact of model grid spacing on regional- and urban-scale air quality predictions of organic aerosol.  *Atm. Chem. Phys.*, 11, 3,107-3,118, 2011.

Valari, M. and Menut, L., Does an increase in air quality models' resolution bring surface ozone concentrations closer to reality?. J. Atm. Ocean. Tech., 25, 1955-1968, 2008.

Vardoulakis, S., Fisher, B. E. A., Pericleous, K.,  Gonzalez-Flesca, N., Modelling air quality in street canyons: a review. *Atm. Env.*, 37, 155-182, 2003.

Wolke, R., Schröder, W., Schrödner, R.,  Renner, E., Influence of grid resolution and meteorological forcing on simulated European air quality: a sensitivity study with the modeling system COSMO--MUSCAT. *Atm. Env.*, 53, 110-130, 2012.

Zhang, J, Moran, M.D., Zheng, Q., Makar, P.A., Baratzadeh, P., Marson, G., Liu, P., Li, S.-M., Emissions preparation and analysis for multiscale air quality modeling over the Athabasca Oil Sands Region of Alberta, Canada, *Atm. Chem. Phys.*, 18, 10459–10481, 2018.

## 6. Appendix A: Model Evaluation Statistics

Table A1. Model comparison statistics

| Metric and Formula | Range | Ideal Score |
|---|---|---|
| $IOA = \begin{cases} 1 - \dfrac{\sum\|M_i - O_i\|}{2(O_i - \bar{O})}, when \sum\|M_i - O_i\| \leq 2(O_i - \bar{O}) \\ \dfrac{2(O_i - \bar{O})}{\sum\|M_i - O_i\|} - 1, when \sum\|M_i - O_i\| > 2(O_i - \bar{O}) \end{cases}$ | [-1,1] | 1 |
| $CoE = 1 - \dfrac{\sum\|M_i - O_i\|}{(O_i - \bar{O})}$ | [-∞, 1] | 1 |
| $MB = \dfrac{1}{N}\sum(M_i - O_i) = \bar{M} - \bar{O}$ | | 0 |
| $MGE = \dfrac{1}{N}\sum\|M_i - O_i\|$ | | 0 |
| $NMB = \dfrac{\sum(M_i - O_i)}{\sum O_i} = \left(\dfrac{\bar{M}}{\bar{O}} - 1\right)$ | | 0 |
| $NMGE = \dfrac{\sum\|M_i - O_i\|}{\sum O_i}$ | | |
| $RMSE = \sqrt{\dfrac{1}{N}\sum(M_i - O_i)^2}$ | | |
| $r = \dfrac{\sum(M_i - \bar{M})(O_i - \bar{O})}{\sqrt{\sum(M_i - \bar{M})^2 \sum(O_i - \bar{O})^2}}$ | [-1.1] | 1 |

The limits on the summations were removed for brevity; all are from i = 1 to N where N is the number of observation-model pairs, $M_i$ is the i'th model value, O is the i'th observation value, and $\bar{M}, \bar{O}$ are the model and observed mean values, respectively.

## 7. Appendix B: Day Versus Night model performance for the different testing methodologies

Table B1. Surface SO$_2$ observations to model comparison, daytime (9:00-18:00) (ppbv).

|      | OS2.5km | OS1km | OS1km$^{A9}$ | OS1km$^{B9}$ | OS2.5km$^{B9}$ | OS1km$^{B49}$ |
|------|---------|-------|--------------|--------------|----------------|---------------|
| IoA  | 0.374 | 0.286 | 0.352 | 0.712 | 0.762 | 0.872 |
| r    | 0.295 | 0.215 | 0.307 | 0.701 | 0.742 | 0.903 |
| NGME | 1.739 | 1.982 | 1.798 | 0.799 | 0.660 | 0.356 |
| GME  | 4.201 | 4.788 | 4.343 | 1.931 | 1.595 | 0.860 |
| CoE  | -0.253 | -0.428 | -0.295 | 0.424 | 0.524 | 0.744 |
| RMSE | 9.317 | 13.388 | 10.275 | 5.171 | 4.652 | 2.996 |
| NMB  | 0.730 | 0.990 | 0.871 | 0.054 | -0.166 | -0.118 |
| MB   | 1.764 | 2.391 | 2.104 | 0.132 | -0.401 | -0.286 |

- 2119 Samples used

Table B2. Surface SO$_2$ observations to model comparison, nighttime (18:00-9:00) (ppbv).

|      | OS2.5km | OS1km | OS1km$^{A9}$ | OS1km$^{B9}$ | OS2.5km$^{B9}$ | OS1km$^{B49}$ |
|------|---------|-------|--------------|--------------|----------------|---------------|
| IoA  | -0.215 | -0.248 | -0.233 | 0.231 | 0.473 | 0.609 |
| r    | 0.204 | 0.206 | 0.205 | 0.339 | 0.421 | 0.620 |
| NGME | 3.143 | 3.281 | 3.215 | 1.896 | 1.300 | 0.964 |
| GME  | 2.061 | 2.152 | 2.108 | 1.243 | 0.852 | 0.632 |
| CoE  | -1.549 | -1.607 | -1.607 | -0.537 | -0.054 | 0.218 |
| RMSE | 5.055 | 5.450 | 5.450 | 3.802 | 2.858 | 2.313 |
| NMB  | 2.166 | 2.328 | 2.328 | 1.076 | 0.394 | 0.361 |
| MB   | 1.421 | 1.527 | 1.527 | 0.706 | 0.258 | 0.230 |

- 3347 Samples used

Table B3. Surface NO$_x$ observations to model comparison, daytime (9:00-18:00) (ppbv).

|      | OS2.5km | OS1km | OS1km$^{A9}$ | OS1km$^{B9}$ | OS2.5km$^{B9}$ | OS1km$^{B49}$ |
|------|---------|-------|--------------|--------------|----------------|---------------|
| IoA  | 0.485 | 0.440 | 0.465 | 0.639 | 0.712 | 0.789 |
| r    | 0.254 | 0.259 | 0.270 | 0.427 | 0.507 | 0.680 |
| NGME | 0.927 | 1.009 | 0.962 | 0.650 | 0.519 | 0.380 |
| GME  | 7.502 | 8.160 | 7.786 | 5.259 | 4.198 | 3.077 |
| CoE  | -0.030 | -0.120 | -0.069 | 0.278 | 0.424 | 0.577 |
| RMSE | 14.843 | 15.811 | 15.571 | 11.272 | 9.982 | 7.964 |
| NMB  | -0.205 | -0.069 | -0.135 | -0.258 | -0.258 | -0.216 |
| MB   | -1.659 | -0.559 | -1.091 | -2.089 | -2.091 | -1.744 |

- 1252 Samples used

Atmospheric
Chemistry
and Physics
Discussions
Open Access

Table B4. Surface NO$_x$ observations to model comparison, nighttime (18:00-9:00) (ppbv).

|      | OS2.5km | OS1km | OS1km[A9] | OS1km[B9] | OS2.5km[B9] | OS1km[B49] |
|------|---------|-------|-----------|-----------|-------------|------------|
| IoA  | -0.016  | -0.050 | -0.045   | 0.275     | 0.511       | 0.587      |
| R    | 0.113   | 0.081 | 0.083     | 0.118     | 0.240       | 0.295      |
| NGME | 1.913   | 1.982 | 1.971     | 1.366     | 0.920       | 0.777      |
| GME  | 17.235  | 17.858 | 17.756   | 12.306    | 8.291       | 7.004      |
| CoE  | -1.032  | -1.105 | -1.093   | -0.451    | 0.023       | 0.174      |
| RMSE | 35.003  | 44.669 | 43.972   | 32.797    | 18.475      | 16.875     |
| NMB  | 0.958   | 0.988 | 0.990     | 0.458     | 0.126       | 0.039      |
| MB   | 8.634   | 8.899 | 8.915     | 4.124     | 1.139       | 0.350      |

- 1862 Samples used

Table B5. Surface O$_3$ observations to model comparison, daytime (9:00-18:00) (ppbv).

|      | OS2.5km | OS1km | OS1km[A9] | OS1km[B9] | OS2.5km[B9] | OS1km[B49] |
|------|---------|-------|-----------|-----------|-------------|------------|
| IoA  | 0.141   | 0.192 | 0.184     | 0.338     | 0.396       | 0.529      |
| r    | 0.166   | 0.215 | 0.211     | 0.327     | 0.367       | 0.504      |
| NGME | 0.660   | 0.621 | 0.627     | 0.508     | 0.464       | 0.361      |
| GME  | 14.427  | 13.568 | 13.703   | 11.111    | 10.143      | 7.901      |
| CoE  | -0.718  | -0.616 | -0.632   | -0.323    | -0.208      | 0.059      |
| RMSE | 21.209  | 20.063 | 20.035   | 16.714    | 15.140      | 12.466     |
| NMB  | 0.587   | 0.542 | 0.557     | 0.454     | 0.414       | 0.326      |
| MB   | 12.839  | 11.854 | 12.187   | 9.918     | 9.050       | 7.121      |

- 864 Samples used

Table B6. Surface O$_3$ observations to model comparison, nighttime (18:00 to 9:00) (ppbv).

|      | OS2.5km | OS1km | OS1km[A9] | OS1km[B9] | OS2.5km[B9] | OS1km[B49] |
|------|---------|-------|-----------|-----------|-------------|------------|
| IoA  | 0.451   | 0.398 | 0.399     | 0.534     | 0.719       | 0.727      |
| r    | 0.526   | 0.541 | 0.557     | 0.642     | 0.784       | 0.784      |
| NGME | 0.706   | 0.775 | 0.773     | 0.600     | 0.361       | 0.352      |
| GME  | 8.326   | 9.132 | 9.116     | 7.070     | 4.258       | 4.145      |
| CoE  | -0.097  | -0.203 | -0.201   | 0.068     | 0.439       | 0.454      |
| RMSE | 11.236  | 12.029 | 11.974   | 10.297    | 6.935       | 7.137      |
| NMB  | 0.492   | 0.624 | 0.651     | 0.510     | 0.262       | 0.296      |
| MB   | 5.799   | 7.359 | 7.668     | 6.008     | 3.088       | 3.491      |

- 1247 Samples used

Atmospheric
Chemistry
and Physics
Discussions

Table B7. Surface PM$_{2.5}$ observations to model comparison, daytime (9:00-18:00) ($\mu$g m$^{-3}$).

|  | OS2.5km | OS1km | OS1km$^{A9}$ | OS1km$^{B9}$ | OS2.5km$^{B9}$ | OS1km$^{B49}$ |
|---|---|---|---|---|---|---|
| IoA | 0.372 | 0.356 | 0.364 | 0.495 | 0.555 | 0.625 |
| r | 0.232 | 0.244 | 0.245 | 0.350 | 0.387 | 0.493 |
| NGME | 0.816 | 0.837 | 0.827 | 0.657 | 0.579 | 0.487 |
| GME | 5.470 | 5.608 | 5.542 | 4.402 | 3.879 | 3.266 |
| CoE | -0.256 | -0.288 | -0.272 | -0.011 | 0.109 | 0.250 |
| RMSE | 9.607 | 10.312 | 10.034 | 8.059 | 7.286 | 6.626 |
| NMB | -0.189 | -0.152 | -0.166 | -0.231 | -0.281 | -0.258 |
| MB | -1.264 | -1.016 | -1.109 | -1.546 | -1.881 | -1.726 |

- 1862 Samples used

Table B8. Surface PM$_{2.5}$ observations to model comparison, nighttime (18:00 to 9:00) ($\mu$g m$^{-3}$)

|  | OS2.5km | OS1km | OS1km$^{A9}$ | OS1km$^{B9}$ | OS2.5km$^{B9}$ | OS1km$^{B49}$ |
|---|---|---|---|---|---|---|
| IoA | 0.193 | 0.170 | 0.173 | 0.337 | 0.471 | 0.528 |
| r | 0.163 | 0.183 | 0.178 | 0.277 | 0.368 | 0.442 |
| NGME | 0.782 | 0.804 | 0.801 | 0.642 | 0.512 | 0.457 |
| GME | 5.313 | 5.466 | 5.444 | 4.367 | 3.483 | 3.105 |
| CoE | -0.614 | -0.660 | -0.653 | -0.326 | -0.058 | 0.057 |
| RMSE | 7.467 | 7.841 | 7.834 | 6.542 | 5.373 | 5.032 |
| NMB | -0.293 | -0.302 | -0.293 | -0.309 | -0.293 | -0.294 |
| MB | -1.992 | -2.050 | -1.989 | -2.098 | -1.991 | -1.995 |

- Samples used