

1 **Response to Referee's comments for "An Evaluation of the Efficacy of Very High Resolution Air-Quality**
2 **Modelling over the Athabasca Oil Sands Region, Alberta, Canada", Matthew Russell et al, submitted to ACP.**
3 **March 6, 2019.**

4 Referee's comments are given in *italic font*, our responses in regular font.

5 **Anonymous Referee 1:**

6 *This paper carried out a series of model-to-model experiments to evaluate the model performances of simulations at*
7 *2.5km and 1km grid-cell, respectively. The comparisons and evaluations between the simulations at 2.5km and 1km*
8 *grid-cell with surface and aircraft observations were also presented. The simulation results indicate that the 1km*
9 *model capture simulation values closer to observations than the 2.5km model, while, in general, the 1km simulation*
10 *has worse performance than the 2.5km simulation. The study is a valuable work, and the design of the simulation*
11 *experiment is rigorous and reasonable. It suggests that using a smaller grid cell size to get better model results is*
12 *unnecessary, since it won't always lead to simulation improvements, while other factors such as meteorology and*
13 *emission should be taken into account to improve the model simulation performance. The paper is well written and*
14 *clearly structured.*

15
16 We thank the reviewer for these positive comments.
17 *The paper could be further improved if the following comments were addressed:*

18
19 1. *Introduction: when talking about the previous research results, most of the cited papers were published*
20 *before 2013. Adding some latest research findings might be better and background for other models such as*
21 *CMAQ model (1. Eder, B., and S. Yu, 2006. A performance evaluation of the 2004 release of Models-3 CMAQ.*
22 *Atmospheric Environment, 40: 4811-4824; 2. S. Yu et al., 2014. Aerosol indirect effect on the grid-scale*
23 *clouds in the two-way coupled WRF-CMAQ: model description, development, evaluation and regional*
24 *analysis. Atmos. Chem. Phys. 14, 11247–11285, doi:10.5194/acp-14-1-2014).*

25
26 We have added a reference to a meta-analysis of a large number of air-quality models (Emery et al., 2017) in the
27 revised manuscript, since the latter deals in a broader sense with model performance evaluation. Eder and Yu 2006
28 does not deal with the issue at hand in our paper; the impact of model resolution on model performance, and
29 therefore should not be referenced here. We have included the reference to Yu et al 2014 in the revised
30 Introduction, however, since it deals with two resolutions (12km and 4km), although more in the sense of the
31 different microphysics parameterizations that must be employed going between those two scales, rather than the
32 impact of resolution within a single scale for which the same microphysics is being used, as carried out in our work.

33
34 2. *Page 7, 2.1.GEM-MACH: The authors need give a detailed description about how "the the physics module of*
35 *the Global Environmental Multiscale meteorological forecast model" was used for different horizontal grid*
36 *cell size. Specially, for grid cell size >4 km, you will need both resolve cloud and subgrid cloud scheme, while*
37 *for grid cell size <4 km such as 2.5 km or 1 km, you will need turn off the subgrid cloud scheme in your*
38 *simulations. Otherwise, you will double count the subgrid cloud effect. This is one of the biggest problems*
39 *for this study. You must clearly describe these physics module in your paper. This may help you to*
40 *understand your results.*

41
42 We have modified the text to make it clear that the two resolutions compared (2.5km and 1km) both use *only* the
43 Milbrandt-Yau (MY) double moment scheme for the cloud microphysics. We also discuss the cloud
44 parameterizations used at the larger scales (note that these are not cross-compared in the paper). The MY scheme
45 was also used in the intermediate 10km resolution simulation before the explicit microphysics scale was reached –
46 this reduces the spin-up time required for the hydrometeors in the explicit microphysics scheme employed at both
47 resolutions compared in the paper. The following text was added/modified:

48
49 *"Four levels of nesting have been employed in our simulations, shown in Figure 2(a). This version of GEM-MACH*
50 *operates on a rotated latitude-longitude coordinate system wherein the position of the coordinate system poles*

51 are set by the user, allowing rotations of the grid with decreasing grid cell size during nesting. The outermost
52 nested grid corresponds to the westernmost two-thirds of the operational GEM-MACH forecasting domain, with a
53 10km grid cell size, and employ a combination of the Kain-Fritsch sub-gridscale convective cloud scheme (Kain and
54 Fritsch, 1990; Kain, 2004) and a Sunqvist (1988) for cloud parameterizations. Within that outer grid is nested a 10
55 km grid cell size western Canada domain (yellow region, Figure 2(a)) which has been rotated to match the
56 horizontal orientation of the Rocky Mountains, and which makes use of a double-moment microphysics scheme
57 (Milbrandt and Yau, 2005a,b) in place of the Sundqvist (1988) parameterization. and which makes use of a double-
58 moment microphysics scheme (Milbrandt and Yau, 2005a,b) in place of the Sundqvist (1988) parameterization. The
59 intention of this intermediate local 10km simulation domain was to provide initial hydrometeors for the two
60 innermost domains, to reduce the “spin-up” time required for the inner domains’ meteorology to reach an
61 equilibrium with respect to cloud formation. The latter two domains (2.5km and 1km grid cell sizes) resolve the
62 cloud microphysics explicitly using the double moment scheme alone and no convective parameterization
63 (Milbrandt and Yau, 2005a,b). ”

64
65 So both 2.5km and 1km resolution are using the same double-moment scheme for their simulations; there is no
66 difference in the cloud microphysics approach used for the model grid cell sizes compared in the paper.

67
68 3. P8, L200-201: *“The forecasts run in a repeating cycle from new meteorological analyses on every 36 hours,
69 and hence are constrained by observations: ”. How can the forecasts be constrained by the observations? I
70 believe that you are doing retrospective simulations other than forecasts. Please explain it clearly.*

71
72 We have clarified the manuscript on this point; the intent was to point out how meteorological analyses with data
73 assimilation underlie the meteorological initial and boundary conditions in our simulations. There are two main
74 considerations: (1) The model simulations were carried out using a configuration of models which mimics an
75 operational forecasting production. (2) In the latter, meteorological observations are utilized in each successive
76 forecast cycle to improve the model initial meteorological state for any given forecast. These same data-
77 assimilated “analysis” fields are the starting points for our model cycling strategy. The use of data assimilation in
78 weather forecasting has been in place in weather forecasting for many years. We have modified our description as
79 follows:

80
81 “Model simulations mimic an operational forecasting system, starting from the use of archived, data-assimilated
82 meteorological analyses as meteorological input and boundary conditions every 36 hours. The use of analysis
83 fields is a standard meteorological forecasting practice to prevent the chaotic drift of the model results from
84 observed meteorology over time. The outermost 10km domain uses initial and boundary conditions from the
85 output of a meteorological simulation, that is itself driven by an analysis field. The outermost domain model then
86 carries out a 36-hour forecast, of which the first 6 hours are discarded as spin-up; the final 30 hours are used as
87 initial and boundary conditions for the rotated 10 km grid cell size domain (the OS10km domain). An OS10km
88 simulation of 30 hours is then carried out...”

89
90 4. P14, regarding the model evaluation metrics, you can refer to the following paper for the definitions these
91 metrics (S. Yu et al., 2006. New unbiased symmetric metrics for evaluation of air quality models.
92 Atmospheric Science Letter, 7, 26-34).

93
94 We had already provided the reference for the metrics in our original manuscript – the “openair” package of
95 Carslaw and Ropkins, 2012. We have added an additional sentence at the same point in the text, “Further
96 discussion of different metrics for model evaluation may also be found in Yu et al., (2006).” The Index of
97 Agreement described in Carslaw and Ropkins is apparently another type of symmetric metric as described in Yu et
98 al. (2006)

99
100 5. Page 16, line 394: “compare” => “comparing”

101 Done.

102

103 6. *Page 17, line 418: "Selecting of " => "Selecting"*

104 The original in the manuscript was "Selection of", not "Selecting of"; the original will be used.

105
106 7. *Page 20, line 515-517: "The expected advantages of the small grid cell size, such as better representation of*
107 *the concentrations of species within plumes and hence better representation of their reactive chemistry (c.f.*
108 *Lonsdale et al., 2012), may be lost in a standard performance analysis due to these other issues.". It*
109 *proposes an important question, is it scientific to evaluate the model performances using the standard*
110 *statistical metrics? Please have more explanations.*

111
112 We would not go so far as to characterize the standard performance metrics as in some way *less scientific*, rather,
113 they specifically evaluate a point-to-point, obs to model comparison. We feel that this is still the most rigorous way
114 to evaluate an air-quality model for performance at specific monitoring locations. However, what we point out in
115 our work is that the additional performance metrics such as those presented in our work show:

- 116 (A) That higher resolution model simulations, even if their standard metric performance has not improved, or
117 is even degraded relative to lower resolution models, nevertheless have scientific value. We have shown
118 that the plume magnitudes are better captured in the high resolution models, if not their location.
119 (B) That the lack of improvement for the standard metrics when going to higher resolution is likely due to
120 difficulties with predicting very localized transport conditions. This is actually a crucial point, in that it
121 suggests that the pathway to improving model performance at higher resolutions should at least in part
122 focus on improving the meteorological transport representation, if needs be through the use of data
123 assimilation of local meteorological observations.

124
125 We have emphasized these points in the conclusions with a short additional paragraph as follows: "Our results
126 should not be taken as an indication that the standard metrics for model comparison are in some way flawed –
127 they provide the most rigorous method for evaluating the performance of a model at specific monitoring locations
128 and specific times. However, the ancillary performance assessment methodology presented here shows that
129 models with very small grid sizes, which may have standard performance metric scores that have not improved or
130 even have degraded relative to larger grid cell size models, nevertheless have scientific value, in terms of being
131 better able to capture plume concentrations and hence plume chemistry, if not plume position. The work also
132 suggests that the prediction accuracy of very local transport conditions may be a large factor in preventing the
smaller grid cell size models from achieving improved performance in standard performance analyses."

133
134 8. *Page 23, line 543-547: "We also noted substantial differences in the day and night performance of both*
135 *models across the methodologies." How to explain the differences?*

136 We have added the following sentences: "The study area is located in a broad river valley with frequent slope-
137 defined anabatic/akatabic and drainage flow events. These often have a diurnal nature, and may explain part of
138 the day/night differences. Example sources of these differences may include the relative ability of the driving
139 meteorological model to capture daytime versus nighttime mixed layer turbulence and planetary boundary layer
140 height."

141
142 9. *Table 2 to table 7, and tables in appendix: "GME" => "MGE" and "NGME"=> "NMGE". Corrected.*

143
144 10. *Page 28, line 613: "Decrease to"=> "Decreasing". Done.*

145
146 11. *Page 28, line 620: "this"=> "the"? We feel that "Despite this difficulty" is better in this context; the*
147 *paragraph thus references the one immediately before it.*

148
149 12. *Acknowledgements, line 635: "with"=> "wish" or "want". Done (used "wish").*

150
151 **Anonymous Referee 2:**

152
153 *The authors present an interesting and well-planned study to evaluate the impacts of high-resolution modeling on*

154 air quality model performance. As growing computational resources facilitate model runs at higher grid resolutions,
155 it is important to understand the extent of the improvements that can be expected from increased resolution and
156 limitations that will continue to constrain model performance, especially if tied to traditional assessment metrics.
157 For this reason, the study conveys a valuable message that should be shared with the modeling community. The
158 study is carefully structured and the manuscript is well-written.

159
160 We thank the reviewer for these positive comments.

161
162 However, several modifications can be made to strengthen the manuscript. Some comments are included below:

- 163
164 1. While the manuscript compares model performances under 2.5 and 1 km grid resolutions, the authors
165 should also discuss model performance relative to acceptable performance benchmarks for air quality
166 modeling. Do the simulations, with either grid resolution, meet recommended performance benchmarks, for
167 example those reported in Emery, et al., 2017 (doi:10.1080/10962247.2016.1265027)? Showing that the
168 modeling was able to meet standard performance expectations would add confidence to the conclusions
169 drawn about the effect of increasing resolution by indicating that the case is an adequate one to draw
170 conclusions from.

171
172 We appreciate the importance and utility of carrying out a meta-analysis of multiple air-quality studies in order to
173 report the current “state of the science” using common air-quality metrics, and of providing a relative ranking of
174 model results at the time of reporting, as was done in Emery et al. (2017). We have added the above reference to
175 our Introduction with that thought in mind. However, we respectfully disagree with the reviewer’s contention
176 that the ability (or failure) for a given modelling study to fall within the percentiles reported in Emery et al. (2017)
177 impacts the confidence of conclusions drawn in a diagnostic study of relative model performance with respect to a
178 specific parameterization or algorithm employed. Emery et al. (2017) state that “The purpose of benchmarks is to
179 understand how good or poor the results are relative to historical model applications of similar nature, and to
180 guide model performance improvements prior to using the model results for policy assessments. To that end, it
181 also remains critical to evaluate all aspects of the model via diagnostic and dynamic methods.” Our study is an
182 example of Emery et al.’s definition of a “diagnostic” evaluation, “in which chemical and physical processes within
183 the model are analyzed individually and collectively”. The work of Emery et al. (2017) provides a useful evaluation
184 of the range of model performance for specific chemical species at the time of writing (specifically, hourly or
185 maximum daily average 8 hour O₃, 24 hour average PM_{2.5} and speciated PM_{2.5}). Their results are useful in the
186 context of operational model evaluation, but they do not add or reduce confidence in diagnostic evaluations such
187 as carried out in our work. Rather, the two approaches are separate, albeit complementary, avenues to identify
188 model performance issues. The comparison of a model score to other models (in other domains, with other
189 emissions data, etc.) allows a ranking relative to past performance. It will not, if that performance is either
190 improved or degraded relative to the historical simulations, necessarily explain the reasons *why* this may be the
191 case, unless a diagnostic evaluation, such as the one we have carried out, is employed. Diagnostic evaluations
192 such as our model grid cell size comparison, however, provide that guidance. With this in mind, we have added
193 the following sentences to the Introduction, “The current state of model science is typically evaluated through
194 multi-model intercomparisons (e.g. Im et al., 2015), and the meta-analysis of these studies can be used to provide
195 useful benchmarks to assess current model performance for specific model species and observations (Emery et al.,
196 2017). However, such studies do not identify the causes for good or poor performance relative to the benchmarks
197 – diagnostic studies, “in which chemical and physical processes within the model are analyzed individually and
198 collectively” (Emery et al., 2017) are required for this purpose. Examinations of the impact of model grid cell size
199 on performance are an example of such a diagnostic evaluation.”

- 200
201 2. The manuscript describes 1km grid simulation as “very high resolution”. However, recent work with
202 regional-scale models such as CMAQ or WRF-Chem has been carried out at horizontal grid resolutions of 1
203 to 3 km. Many of the modeling studies referenced in the manuscript are several years old. A deeper
204 discussion of the progression and current state of grid resolution in Eulerian air quality modeling would

strengthen the paper. The paper should discuss what constitutes “very high resolution” at present and, more importantly, what maximum level of resolution can be expected from existing modeling frameworks given the dependence of existing subgrid-scale parameterizations on grid resolution.

We have added the following text to the Introduction:

“Air-quality model grid-cell size typically follows the grid-cell sizes used in weather forecasting models, which have followed a gradual progression towards finer discretization where more explicit representation of cloud formation and local radiative transfer effects may be represented. The most recent weather forecasting applications (e.g. Leroyer *et al.*, 2014) have reached grid-cell sizes as small as 250m over limited domains such as individual cities, and have shown promising results in terms of being able to resolve some aspects of local circulation. In addition, as grid resolution reaches the 3 to 4 km scale, explicit cloud microphysics packages may be used, allowing potentially better performance, particularly with regards to feedbacks between meteorology and chemistry (Yu *et al.*, 2014; Gong *et al.*, 2015). However, while these models promise better physical representation of local chemistry, their performance may be limited by the quantity and availability of initialization and boundary condition meteorological data; these data may be used in a data assimilation context to improve their initial state. The accuracy of broader-scale meteorological predictions may thus influence local model accuracy, despite the ongoing reduction in meteorological model (and consequently air-quality model) grid cell size. Some recent air-quality model simulation studies with grid cell sizes on the order of one to four km include Thompson and Selin (2012), Li *et al.* (2014), Joe *et al.* (2014), Kheirbek *et al.* (2014), Kheirbek *et al.* (2016), and Pan *et al.*, (2017).

3. Although the manuscript’s analysis is well structured, some additional discussion of how the findings can be expected to be representative of air quality modeling beyond this specific simulation would be beneficial. Do the authors expect the findings to remain consistent across often applied increasing resolution levels in regional-scale air quality modeling, for example 36km to 12km to 4 km? Should similar conclusions be expected over more urban domains?

This is a very good question, though difficult to answer quantitatively without broadening the scope of our original study significantly. We have added the following paragraph as a short Discussion section before the Conclusions, to try to address the issue in qualitative sense:

“A key result of our current work is that 1km grid cell size simulations resulted in improved prediction of plume concentration maxima relative to 2.5km grid cell size simulations, despite having no improvement using standard scoring methodologies. We also have described a scoring approach wherein these potential advantages of higher resolution may be quantified. We believe that flow field effects such as described in Figure 3 are a general result of increasing grid resolution, but note important caveats, which include:

- (1) The availability of meteorological observation and high resolution emissions data to provide model driving information, and the resolution and proximity of this information to the simulation location. Both will influence the relative importance of grid cell size on model results. If this information is available in a higher resolution than the lower of two grid cell size simulations being compared, and/or is used via data assimilation to improve model initial meteorological conditions, our expectation is that the smaller grid cell size model may outscore the larger grid cell size model, even for more standard metrics.
- (2) The extent to which local, versus synoptic, weather conditions drive flow in a given region. For example, in the urban heat island meteorological simulations of Leroyer *et al.* (2014), the accuracy of local flow predictions was shown to be extremely dependent on the representation of the urban heat island, and the accuracy of the latter was critically dependent on the grid cell size (which in this example went down to 250 m). In this respect, for meteorological conditions wherein local factors can dominate the flow, and where those conditions may be adequately modelled only at very high resolution, we would again expect the smaller grid cell size simulation to provide better performance, for standard metrics.
- (3) Conversely, model performance using standard metrics should not be expected to *increase* with successively larger and larger grid sizes; the accuracy of even the synoptic flow field will not be captured as model resolution decreases.

Given these considerations, we recommend that modellers should attempt successively smaller grid cell sizes to

257 determine the following: first, the point at which, for their particular system and simulation location, subsequent
258 grid cell size reductions fail to improve performance; and second, to make use of still higher resolutions for studies
259 wherein the point-to-point comparison is less important, and other factors such as accurately capturing the plume
260 chemistry are more crucial.”

- 261
- 262 4. *The manuscript states that the study results are “strongly suggestive of the presence of issues such as
263 illustrated in Figure 3”, that is plume structures that are better represented by the higher resolution but
264 more affected by errors in wind fields. An illustrative example of this taken from the simulated results would
265 strengthen this conclusion. A comparison of simulated plumes that mirrors the schematic included in figure
266 3 would be beneficial.*

267 We have modified Figure 8 in the manuscript, since it provides such an example, and added the following text to
268 the description of Figure 8 in the revised manuscript: “Panels (a) and (c) of Figure 8 provide a further example of
269 the kind of situation referenced in Figure 3; surface monitoring station locations are depicted as grey circles, one of
270 which is identified with a pink arrow. This station lies within the plume at 2.5km resolution (Figure 8(a)), and
271 outside of the plume at 1km resolution (Figure 8(c)). While the plume direction is the same at both scales, that is,
272 the large-scale wind field controls the positioning of the plume axis, the smaller grid cell size simulation places a
273 stronger constraint on the accuracy of the wind field. For example, if the simulated large-scale flow direction was
274 inaccurately predicted by only a few degrees, the plume would not appear in the 1km simulation time series at this
275 location, while registering as present in the 2.5km simulation. Nevertheless, the plume maximum concentration is
276 better captured by the smaller grid cell size simulation (compare maximum values in observed aircraft SO₂, Figure 8
277 (b, d)). The higher resolution simulation may thus more accurately simulate the plume maximum concentration –
278 but not its placement in space, as was hypothesized in Figure 3.”

- 280
- 281 5. *The authors briefly mention the connection between grid resolution in air quality modeling and associated
282 health impacts projections (line 67-70). Previous work has looked at the impact of increasing grid resolution
283 and improved model performance on health effects estimates, and how these sources of uncertainty
284 compare (e.g., Thompson, et al., 2012, doi:10.5194/acp-12-9753-2012). Some additional discussion of the
285 role of uncertainty due to grid resolution in the larger context of air quality impact assessments, including
286 exposure and health impacts, would be beneficial.*

287 We thank the reviewer for that reference, which has been added to the revised manuscript. Thompson and Selin
288 (2012) noted that they found that increases in resolution to grid cell sizes below 12 km did not improve their
289 health outcome predictions. This may have been due to the same issues as we have noted in our work: higher
290 spatial resolution does not necessarily guarantee a more accurate prediction at the locations at the observation
291 locations. A lower correlation score at decreasing grid size for example would imply a greater degree of difficulty
292 with linking model output to health effects. However, the studies quoted noted that *accurate* higher resolution
293 simulations are nevertheless *desired* for health studies, due to the need to relate exposure to concentrations on
294 the scale of a few kilometers. Our work helps to explain Thompson and Selin (2012)'s findings, and it points to a
295 possible way to further improve high-resolution model results, through local data assimilation, as we noted in our
296 conclusions. We have modified the given sentence in the Introduction to include this information:

298

299 “These studies have often demonstrated that failure to account for higher resolution features may result in
300 mischaracterization of concentrations or health impacts (Isakov et al., 2007), although the capability of current
301 models to provide this information with sufficient accuracy is unclear. One study found that increasing resolution
302 did not change predicted health outcomes, and concluded that “resolution requirements should be assessed on a
303 case-by-case basis” (Thompson and Selin, 2012), while others (e.g. Kheirbek et al. (2014), Kheirbek et al. (2016))
304 have employed 1km resolution without discussing the impacts of resolution on predicted health outcomes.”.

306 We have also modified extended the paragraph to include a new closing sentence:
307
308 "The health studies carried out to date highlight the need for better understanding the underlying controlling
309 factors for model accuracy with decreasing grid cell size."
310
311 *Smaller comments:*
312 - Lines 64-65: This sentence is unclear.
313
314 The sentence has been modified as follows: "A number of studies have tried to evaluate the benefits of higher
315 resolution simulations and to quantify the impact of sub-grid variability by using different model grid-cell sizes".
316
317 - Lines 97-99: Expand on this statement. What specifically makes the VHR representations more realistic?
318 We have added the following lines of text:
319 "Salvador *et al.* (1999) studied the prediction accuracy impacts of meteorological model grid cell size in a region
320 with complex domain, and found that 2km or smaller grid cell sizes were required to resolve local scale complex
321 terrain flow features, and that daytime vertical advection and predictions of turbulent kinetic energy and potential
322 temperature were influenced by grid cell size. Dore *et al.* (2012) evaluated air quality model NO₂ simulations
323 employing 1, 5 and 50km grid cell sizes against observations, and found the best performance for the 1km
324 simulation, with more physically realistic distributions of reactive nitrogen, attributing this performance gain to
325 more realistically precipitation simulations and emissions inputs for the smallest grid cell size. The availability of
326 high-resolution emissions information may be a limiting factor in improved simulations as grid cell size decreases.
327 Valari and Menut (2008) noted that emissions inaccuracy was the principal cause of noise in small grid cell size
328 simulations conducted for the Paris area, and proposed the use of statistical downscaling in favour of predictive
329 modelling at scales at or below 1km grid cell size."
330 - Line 231: Remove "for areas"
331 Done.
332 - Line 462: Changing "first three columns" to "third column", might be clearer
333 Done.
334 - Tables 2-5: Including the definition of each acronym used for the metrics somewhere on the chart or at the
335 beginning of the charts would improve readability.
336 It turned out that the word version of the acronym could fit in the tables, so that has been used.
337
338 - Line 567-569: The issue of air quality models excessively mixing pollutants along the vertical dimension within the
339 boundary layer has been previously acknowledged by several studies (e.g. Garcia-Menendez, *et al.*, 2014,
340 doi:10.1016/j.scitotenv.2014.05.108).
341 Thanks, we have added that reference to this discussion: "Garcia-Menendez *et al.* (2014) have noted similar
342 results for forest fire plume prediction."
343
344 - Figure 8 needs to be improved. The x-axis of the left panels is illegible. Lines and colors on the right plots are a bit
345 hard to observe as well; a higher resolution/quality plot would help.
346
347 Several improvements were done to Figure 8 in the revised version – we used a smaller time interval for panels (b)
348 and (d) in order to better show the differences in the simulations, and increased the size of the font for the axes of
349 all panels and colour bars. We have also modified panels (a) and (c) to address show the similarities between the
350 simulated concentration fields at each resolution and the hypothesized impact of flow inaccuracies given in Figure
351 3. In the process of updating this Figure, we also noticed that the 11 UT contour field had been used for panels (a,c)
352 rather than the intended 11 MDT contour field – this was corrected as part of the updates to this Figure. |

353
354 **An Evaluation of the Efficacy of Very High Resolution Air-Quality**
355 **Modelling over the Athabasca Oil Sands Region, Alberta, Canada**

Formatted: Normal, Space Before: 0 pt, Line spacing: single,
Allow hanging punctuation
Formatted: Indent: Left: 0 cm
Formatted: Font: Bold

356 Matthew Russell¹, Amir Hakami¹, Paul A. Makar², Ayodeji Akingunola², Junhua Zhang², Michael D. Moran², and
357 Qiong Zheng²

358 ¹Department of Civil and Environmental Engineering, Carleton University, Ottawa, Canada

359 ²Air Quality Research Division, Environment and Climate Change Canada, Toronto, Canada

360
361 **Abstract**

362 We examine the potential benefits of very high resolution for air-quality forecast simulations using a nested
363 system of the Global Environmental Multiscale – Modelling Air-quality and Chemistry chemical transport model.
364 We focus on simulations at 1km and 2.5km grid-cell spacing for the same time period and domain (the industrial
365 emissions region of the Athabasca Oil Sands). Standard grid-cell to observation station pair analyses show no
366 benefit to the higher resolution simulation (and a degradation of performance for most metrics using this
367 standard form of evaluation). However, when the evaluation methodology is modified, to include a search over
368 equivalent representative regions surrounding the observation locations for the closest fit to the observations, the
369 model simulation with the smaller grid cell size had the better performance. While other sources of model error
370 thus dominate net performance at these two resolutions, obscuring the potential benefits of higher resolution
371 modelling for forecasting purposes, the higher resolution simulation shows promise in terms of better aiding
372 localized chemical analysis of pollutant plumes, through better representation of plume maxima.

373 **1 Introduction**

374 Numerical modeling of the atmosphere in an Eulerian framework relies on discretization of the computational
375 domain into a numerical grid. The horizontal grid cell size of atmospheric simulations can range in from hundreds
376 of kilometers, to the metre-scale of Large Eddy Simulation models. Air-quality model grid-cell size typically
377 follows the grid-cell sizes used in weather forecasting models, which in turn have followed a gradual progression
378 towards finer discretization, where more explicit representation of cloud formation and local radiative transfer
379 effects may be represented. The most recent weather forecasting applications (e.g. Leroyer et al., 2014) have
380 reached grid-cell sizes as small as 250m over limited domains such as individual cities, and have shown promising
381 results in terms of being able to resolve some aspects of local circulation. In addition, as grid resolution reaches
382 the 3 to 4 km scale, explicit cloud microphysics packages may be used, allowing potentially better performance,
383 particularly with regards to feedbacks between meteorology and chemistry (Yu et al., 2014; Gong et al., 2015).
384 However, while these models promise better physical representation of local chemistry, their performance may
385 be limited by the quantity and availability of initialization and boundary condition meteorological data; these data
386 may be used in a data assimilation context to improve their initial state. The accuracy of broader-scale

Formatted: Not Highlight
Formatted: Not Highlight
Formatted: Not Highlight
Formatted: Font: Italic, Not Highlight
Formatted: Not Highlight

Formatted: Font: Italic
Formatted: Font: Italic
Formatted: Not Highlight
Formatted: Not Highlight

387 meteorological predictions may thus influence local model accuracy, despite the ongoing decrease in
388 meteorological model (and consequently air-quality model) grid cell size. Some recent air-quality model
389 simulation studies with grid cell sizes on the order of one to four km include Thompson and Selin (2012), Li *et al.*
390 (2014), Joe *et al.* (2014), Kheirbek *et al.* (2014), Kheirbek *et al.* (2016), and Pan *et al.*, (2017).

391 For the purposes of this study, Very High Resolution (VHR) modelling refers to the current higher resolution limits
392 of chemical transport models (CTMs), employing a horizontal grid cell spacing of 1km or less. It is in this regime
393 that the photochemical processes may be forecasted with resolved microphysics (e.g. Milbrandt and Yau,
394 2005(a,b)), and detailed particle and gas-phase chemistry, using currently available computer technology. VHR
395 modelling is very computationally expensive, and also introduces its own set of challenges, such as the availability
396 of surface boundary condition fields as the model grid cell size decreases. Moreover, it is not currently clear
397 whether decreases in model grid cell size leads to more accurate results when compared to observations. The
398 motivation behind VHR modelling in CTMs is to reduce the impact of diluting chemical concentrations - especially
399 from averaging emission plumes into large grid cells – in order to better capture inhomogeneities in emission
400 profiles, to better simulate local transport processes associated with terrain that would otherwise be smoothed by
401 the use of a coarse grid, and to reduce truncation errors and hence achieve better numerical accuracy (Jacobson,
402 1999).

403 We note here that while the terms “grid cell size” and “resolution” tend to be used interchangeably in the
404 literature, this is not true in a precise mathematical sense; more formally, the ability to resolve features of size
405 $2\Delta x$ requires a grid cell spacing of size Δx , and the highest spatial frequency which can be reconstructed from a
406 discrete sampling of the latter grid cell spacing will be $\frac{1}{2\Delta x}$, the Nyquist wavenumber of the grid cell size
407 discretization. Furthermore, atmospheric models may make use of energy dissipation techniques that broaden
408 the size of resolvable wavelengths to $3\Delta x$ to $4\Delta x$ (Grasso, 2000; Pielke, 2001). Model resolution is thus a function
409 of, but not equivalent to, grid cell size. Here, we define “resolution” as the ability of a model to clearly distinguish
410 components of a predicted atmospheric variable, as a *function* of grid cell size.

411 The issue of a model to distinguish these features is also compounded by uncertainties in model inputs. For
412 example, in a large rural setting, a large model grid cell will represent an area containing many roads, whose
413 emissions will be averaged into one value per species per time. As the grid cell size decreases however, this
414 averaging effect will be reduced, giving each road’s emissions more impact on the resulting concentrations in the
415 grid cell containing it. However, the smaller grid cell size will also result in steeper concentration gradients in the
416 model between adjacent grid cells, which can in turn result in numerical instabilities that contaminate predictions
417 ([Salvador, et al., 1999](#)). At the same time, a reduction in grid-cell size can be shown
418 formally to reduce inaccuracies in the discretization of the governing equations for atmospheric motion (Coiffier,
419 2011). Previous efforts to address these issues through variable grid size or structure in air quality modeling have

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Not Highlight

Formatted: Font: Italic

Formatted: French (France)

Formatted: Font: Italic, French (France)

Formatted: French (France)

420 not received sustained attention, and therefore most current air quality models use a uniform (albeit nested) grid
421 cell size in applications (Garcia-Menendez *et al.*, 2010; Kumar *et al.*, 1997).

422 As resolution increases further, the shape-presence of local topographical features (*e.g.* buildings and street
423 canyons) become more important. Both the increased topographic complexity, and potential numerical
424 instabilities can lead to differences in meteorological forcing as resolution increases (Wolke, *et al.*, 2012; Gego, *et*
425 *al.*, 2005)). The contribution of meteorological uncertainties due to resolution become more significant, especially
426 for secondary pollutants such as ozone (Valari and Menut, 2008) or secondary Particulate Matter (PM). For
427 example, Markakis *et al.* (2015) in their analysis of 4 km CHIMERE simulations for the relatively flat terrain of Paris,
428 France, suggested that model meteorological grid cell size does not significantly impact forecast accuracy. That
429 may not have been the case, had their terrain been more complex. In contrast, Queen and Zhang (2008) observed
430 considerable meteorological sensitivity to the more complex terrain in their 4 km resolution Community
431 Multiscale Air Quality (CMAQ, EPA 1999) model simulations simulation over the Appalachian Mountains in the
432 eastern United States, as did Salvador *et al.* (1999) for meteorological model simulations.

Formatted: Font: Italic

433 A number of studies have tried to evaluate the benefits of higher resolution simulations and to quantify the
434 impact of sub-grid variability by using different model grid-cell sizes A number of studies, employing various
435 approaches, have tried to evaluate the benefits of higher resolution simulations by quantifying sub-grid variability
436 by employing larger model grid cell sizes (Vardoulakis *et al.*, 2003; Ching *et al.*, 2006; Pepe *et al.*, 2016). These
437 studies have often demonstrated that failure to account for higher resolution features may result in
438 mischaracterization of concentrations or health impacts (Isakov *et al.*, 2007), although the capability of current
439 models to provide this information with sufficient accuracy is unclear. One study found that increasing resolution
440 did not change predicted health outcomes, and concluded that “resolution requirements should be assessed on a
441 case-by-case basis” (Thompson and Selin, 2012), while others (*e.g.* Kheirbek *et al.* (2014), Kheirbek *et al.* (2016))
442 have employed 1km resolution without discussing the impacts of resolution on predicted health outcomes These
443 studies have often demonstrated that failure to account for higher resolution features may result in
444 mischaracterization of concentrations or health impacts (Isakov *et al.*, 2007). Population exposure studies using
445 air pollution models may be affected by resolution in a more complex fashion, given that both the predicted field
446 (a pollutant with a known health impact) and the data to which the predicted field is to be linked (the human
447 population) both have resolution dependencies. The health studies carried out to date highlight the need for
448 better understanding the underlying controlling factors for model accuracy with decreasing grid cell size.

449 Terrain and meteorology are not the only factors that contribute to greater uncertainties as horizontal grid cell
450 size is reduced – for example, the ability of the model to locally resolve emission fluxes may also become a factor.
451 This may result in improved or deteriorated model performance as the size of the grid cells decrease. Gridded
452 model emissions may have an intrinsic resolution dependence in the underlying spatial disaggregation fields, and

453 this can contribute to uncertainties and errors in emissions as grid cell size is decreased. For instance, Valari and
454 Menut (2008) found that the discrepancy between their modelled and observed concentrations grew, rather than
455 shrank, in response to decreases in grid cell size from 48km to 6 km, and they associated these results with
456 changes in the resulting local emission fluxes. They showed that in their model setup, with regard to ozone, a grid
457 cell size was reached ($12 \times 12 \text{ km}^2$) where errors in inputs (errors in the emission inventory, wind direction, etc.)
458 outweighed the importance of other sources of model error such as grid cell size. The authors however noted that
459 Paris' ozone photochemistry very often resides on the transition between a NO_x^- sensitive and a VOC-sensitive
460 regime (Sillman et al., 2003). These are chemical conditions which can alternatively produce or titrate ozone, and
461 hence have has a degree of sensitivity to precursor emissions, and therefore, also, to any errors in those emissions.
462 Conversely, in a 3-level nested 9- to 3- to 1- km MM5–CMAQ simulation over Osaka, Japan, Shrestha et al., (2009)
463 found that ozone comparisons to observations improved as the grid resolution increased. This was also the case
464 for a 36- to 12- to 4-km nested MM5–CMAQ simulation over Houston, USA (Ching et al., 2006), where the ozone
465 forecast improvement associated with higher resolution was attributed to the ability of the finer grid cell size
466 model nests to adequately resolve high concentrations of freshly emitted NOx and hence allow for more local
467 ozone titration. The latter process might not take effect until the grid cell size is sufficiently fine to resolve the NOx
468 source patterns (i.e., a level where traffic and industrial sources can be identified.) This titration was not seen until
469 they decreased their grid cell sizes to 2 km and smaller. Stroud et al. (2011) noted a similar grid cell size
470 dependent chemical impact on model performance, where secondary organic aerosol formation maxima were
471 better simulated with a 2.5km grid cell size model than a 10km grid cell size model. In general, the impact of
472 resolution on model performance appears to depend on a number of factors, such as the terrain, spatial
473 distribution of sources, pollutant of concern, season, etc. (Arunachalam et al., 2006; Queen and Zhang, 2008; Dore
474 et al., 2012).

Formatted: French (France)

475 Salvador et al. (1999) studied the prediction accuracy impacts of meteorological model grid cell size in a region
476 with complex domain, and found that 2km or smaller grid cell sizes were required to resolve local scale complex
477 terrain flow features, and that daytime vertical advection and predictions of turbulent kinetic energy and potential
478 temperature were influenced by grid cell size. Dore et al. (2012) evaluated air quality model NO_2 simulations
479 employing 1, 5 and 50km grid cell sizes against observations, and found the best performance for the 1km
480 simulation, with more physically realistic distributions of reactive nitrogen, attributing this performance gain to
481 more realistically precipitation simulations and emissions inputs for the smallest grid cell size. The availability of
482 high-resolution emissions information may be a limiting factor in improved simulations as grid cell size decreases.
483 Valari and Menut (2008) noted that emissions inaccuracy was the principal cause of noise in small grid cell size
484 simulations conducted for the Paris area, and proposed the use of statistical downscaling in favour of predictive
485 modelling at scales at or below 1km grid cell size. Whether or not simulated quantities improve with reference to
486 observed quantities as applications approach VHR grid cell sizes, the resulting distribution of the quantities tends

Formatted: English (United States)

Formatted: Font: Italic, English (United States)

Formatted: English (United States)

Formatted: Font: Italic

Formatted: Not Highlight

Formatted: English (United States)

487 to be more physically realistic (Dore *et al.*, 2012; Salvador *et al.*, 1999; Valari and Menut, 2008).

488 The current state of model science is typically evaluated through multi-model intercomparisons (e.g. Im *et al.*,
489 2015), and the meta-analysis of these studies can be used to provide useful benchmarks to assess current model
490 performance for specific model species and observations (Emery *et al.*, 2017). However, such studies do not
491 identify the causes for good or poor performance relative to the benchmarks – diagnostic studies, “in which
492 chemical and physical processes within the model are analyzed individually and collectively” (Emery *et al.*, 2017)
493 are required for this purpose. Examinations of the impact of model grid cell size on performance are an example
494 of such a diagnostic evaluation.

Formatted: Font: Italic

Formatted: Font: Italic

Formatted: Font: Italic

495 The benefits for model performance with increased spatial resolution are unclear, based on the above literature.
496 However, most papers converge towards the following qualitative conclusions:

- 497 1. The impact of terrain topology on meteorological forcing as grid cell size decreases can dwarf the impact of
498 a more accurate spatial apportionment of the corresponding emissions.
- 499 2. Decreases in grid cell size result in a more realistic spatial distribution of chemical species, whether or not
500 model performance is improved.
- 501 3. Uncertainties of spatial and temporal emissions allocation have an increasing influence on overall model
502 uncertainty as model grid cell size decreases.

503 The 1980's saw several studies in which the potential impacts of wind direction errors on dispersion model
504 performance were examined. Fox (1981) noted that pairing of model output at observation station locations could
505 be done as a function of both time and space *e*, as a function of time *(-by* {combining the data across all stations*)},*
506 as a function of space *(by* {combining all times, at each station location*)*, or without any pairing (observations and
507 data *were* compared as cumulative frequency distributions). The accuracy of regulatory dispersion models in the
508 early 1980's was such that Fox (1984) concluded that model and observation values paired in time and space
509 exhibited “little to no correlation” and discussed potential errors associated with transport. Poor correlations were
510 also noted by Hanha (1988), reporting on the first generation of reactive-transport models, stated “wind direction
511 errors are the major cause of the poor agreement in hourly predictions of concentrations at short distances
512 downwind of point sources,” as well as describing metrics for air-quality model evaluation. Hanha (1988) also
513 noted that model predictions could be offset in space and time relative to observations, leading to poor
514 performance statistics, despite a greater degree of similarity of behavior if the offsets are taken into account.
515 Errors in wind-field modelling were described as the main source of error in simulations of plumes by Carhart *et al*
516 (1989), again showing how better agreement resulted when model and observations were unpaired in time and/or
517 space, and noted that other metrics such as maximum plume width might better represent model performance.

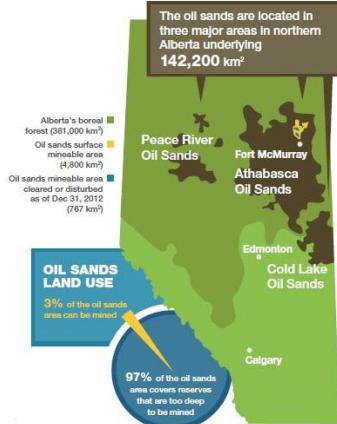
518 Lee (1987) found that small perturbations in space and time could result in poor correlations, despite similar
519 histogram distributions of both model and observations.

520 More recently, Kang *et al.*, (2007) examined the concept of using the area of the limiting resolution of the model (2
521 to $3\Delta x$, where Δx is the horizontal grid cell size) to weight or spatially average model evaluation metrics for a single
522 grid-cell size, noting how the model's rated ability to capture high concentration events ("hits") was increased when
523 the limiting resolution of the model was incorporated into the performance metrics. However, the use of averaging
524 may mask the potential for a model with a small grid cell size to contain both the desired plume magnitude, as well
525 as much lower concentrations, within the same larger representative area, in turn masking the potential impact of
526 the reduction in grid cell size.

527 We expand on this concept to evaluate the impact of model grid cell size in the context of an equivalent area about
528 a given observation location. We examine area-weighted metrics in the form of averages over roughly equivalent
529 areas for different model grid cell sizes, and also use the *a priori* knowledge of the observations to determine
530 whether the closest match to observations may be found within an equivalent area. We show that the latter metric
531 demonstrates a positive impact of model grid cell size on simulation results, while more simple paired comparisons,
532 and averages over similar areas, mask these benefits.

533 We examine the impact of grid cell size on model performance in a region of intense petrochemical extraction and
534 upgrading, the Athabasca Oil Sands Region (AOSR). The AOSR refers to the northernmost of three large bitumen
535 deposits located the northern part of the province of Alberta in Canada; the Athabasca, Peace River, and Cold Lake
536 areas. Together these areas cover 142,200 km² in total, and constitute the third largest oil reserves in the world
537 (Government of Alberta, 2016), as shown in Figure 1. The oil sands sector is the second largest source of SO₂ and
538 the third largest source of industrial NO_x in the province of Alberta. This sector is also a significant source of
539 industrial PM, CO, and Volatile Organic Compound (VOC) emissions (Zhang *et al.*, 2018), from a variety of source
540 types and industrial processes (e.g. open pit mine tailings ponds, large diesel fleets, bitumen upgrading facilities).
541 As is described below, very high resolution emissions data are available for these sources, and emissions take place
542 in a region with significant topography, hence the region provides a good test case for the relative impact of grid
543 cell size on air-quality model prediction results.

544 We describe next our model, the simulation domains and forecasting setup, the emissions data, our evaluation
545 methodology, and the results of our analysis.



546

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

548

2.

← **Formatted:** Indent: Left: 0 cm, Hanging: 0.76 cm, No bullets or numbering

Formatted: Normal

549 Methodology

550 2

551 2.11.1 GEM-MACH

552 The air-quality model used in this work is Environment and Climate Change Canada's (ECCC) Global Environmental
553 Multiscale – Modelling Air-quality and Chemistry (GEM-MACH) model, which has been in use as Canada's
554 operational air-quality forecast model since 2009 (Moran *et al.*, 2010). GEM-MACH is an on-line model, that is,
555 both meteorological and chemistry processes are handled within a single model. The chemical processes reside
556 within the physics module of the Global Environmental Multiscale meteorological forecast model (Côté, *et al.*,
557 1998(a,b)), originate with Environment Canada's earlier off-line model (A Unified Regional Air-quality Modelling
558 System; AURAMS, Gong *et al.*, 2006), and include process representation for particle microphysics (Gong *et al.*,
559 2003(a,b)), inorganic heterogeneous chemistry (Makar *et al.*, 2003), aqueous phase chemistry, in-cloud and below-
560 cloud scavenging (Gong *et al.*, 2006), and secondary organic aerosol formation (Stroud *et al.*, 2011). GEM-MACH
561 employs a sectional approach to represent the size distribution of atmospheric particles, with 12-bin (Makar *et al.*,
562 2015(a,b); Gong *et al.*, 2015) or 2-bin configurations (Moran *et al.*, 2010). The latter configuration is designed for
563 maximum computational efficiency, with re-binning to the 12-bin distribution for key particle microphysics
564 processes, in order to improve accuracy. Here, the 2-bin version of the model has been used, the main focus of the
565 work being the impact of horizontal grid cell size on model results. Eight aerosol chemical components are resolved
566 in GEM-MACH (sulphate, nitrate, ammonium, elemental carbon, primary organic aerosol, secondary organic
567 aerosol, sea-salt and crustal material). In the present study, we make use of GEM-MACH v.1.5.1, described in more
568 detail in Makar *et al.*, 2015(a,b), employing 80 levels in a hybrid vertical coordinate system extending up to 0.1hPa
569 (~30km). Both model grid cell size simulations compared here (2.5km and 1km grid cell sizes, see below) make use
570 of the Milbrandt-Yau double moment explicit microphysics scheme, that is, cloud processes are resolved explicitly
571 at these scales (Milbrandt and Yau, 2005(a,b)).–

572 2.21.2 Model Setup

573 2.21.2.1 Grid Nesting

574 Four levels of nesting have been employed in our simulations, shown in Figure 2(a). This version of GEM-MACH
575 operates on a rotated latitude-longitude coordinate system wherein the position of the coordinate system poles is
576 set by the user, allowing rotations of the grid with decreasing grid cell size during nesting. The outermost nested
577 grid corresponds to the westernmost two-thirds of the operational GEM-MACH forecasting domain, with a 10km
578 grid cell size, and employ a combination of the Kain-Fritsch sub-gridscale convective cloud scheme (Kain and
579 Fritsch, 1990; Kain, 2004) and a Sundqvist (1988) for cloud parameterizations. Within that outer grid is nested a 10
580 km grid cell size western Canada domain (yellow region, Figure 2(a)) which has been rotated to match the
581 horizontal orientation of the Rocky Mountains, and which makes use of a double-moment microphysics scheme
582 (Milbrandt and Yau, 2005a,b) in place of the Sundqvist (1988) parameterization. The intention of this
583 intermediate local 10km simulation domain was to provide initial hydrometeors for the two innermost domains.
584

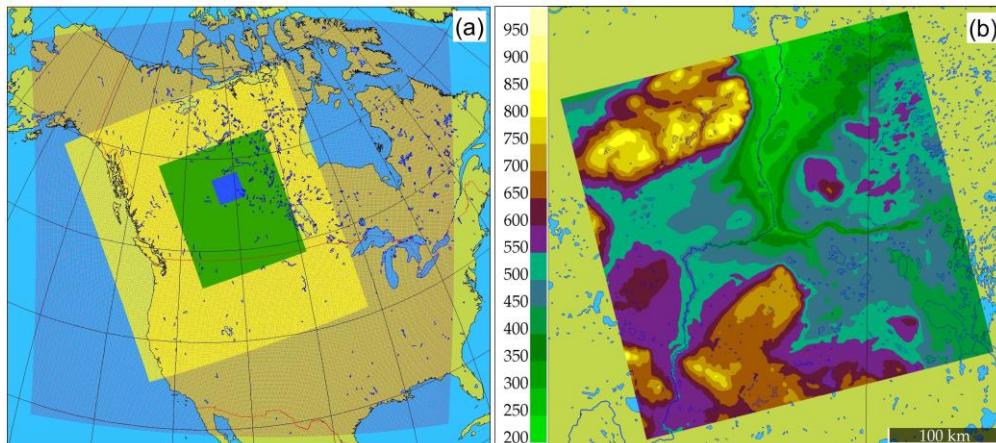
585 to reduce the “spin-up” time required for the inner domains’ meteorology to reach an equilibrium with respect to
586 cloud formation. The latter two domains (2.5km and 1km grid cell sizes) resolve the cloud microphysics explicitly
587 using the double moment scheme alone and no convective parameterization (Milbrandt and Yau, 2005a,b). Four
588 levels of nesting have been employed in our simulations, shown in Figure 2(a). This version of GEM-MACH
589 operates on a rotated latitude longitude coordinate system wherein the position of the coordinate system poles
590 may be set by the user, allowing rotations of the grid with decreasing grid cell size during nesting. The outermost
591 nested grid corresponds to the westernmost $\frac{2}{3}$ of the operational GEM-MACH forecasting domain, with a 10km
592 grid cell size. Within that is nested a 10km grid cell size western Canada domain (yellow region, Figure 2(a)) which
593 has been rotated to match the horizontal orientation of the Rocky Mountains, and which makes use of a similar
594 double moment microphysics scheme (Milbrandt and Yau, 2005 (a,b)) as the two innermost domains. — this
595 intermediate nested grid was constructed in order to allow hydrometeors to be passed from the western Canada
596 10km domain to the two innermost domains with a minimum of spin up time required for the inner domain’s
597 meteorology. The third nested grid inwards (green region, Figure 2(a)) is the 2.5km grid cell size domain, which
598 covers most of the Canadian provinces of Alberta and Saskatchewan. This grid will hereafter be referred to as the
599 OS2.5km domain. The fourth and final nested grid (blue square, Figure 2(a)) is a 1km grid cell size domain, roughly
600 centered over and covering the immediate environs of the Athabasca Oil Sands, and is referred to hereafter as the
601 OS1km model. This last nest also shows the region within which 22 instrumented aircraft flights were conducted
602 during August and September of 2013, providing a unique measurement dataset for our evaluation of the
603 OS2.5km and OS1km model output for the same time period. Table 1 provides details on the horizontal
604 dimensions of each of these nested domains, and the duration of the simulations on each grid. All four model
605 nests make use of the same vertical coordinate and levels. Figure 2(b) shows the topography of the 1km domain
606 in detail; the region to be modelled is situated in a broad river valley, with a local vertical relief of 750 m.
607 Significant wind shears and frequent inversions are observed in the region, and part of our interest in 1km grid cell
608 size simulations is to determine the extent to which these local features may influence model prediction accuracy.

609 2.2.2 Simulation Cycling Strategy

610 Model simulations mimic an operational forecasting system, starting from the use of archived, data-assimilated
611 meteorological analyses as meteorological input and boundary conditions every 36 hours. The use of analysis
612 fields is a standard meteorological forecasting practice to prevent the chaotic drift of the model results from
613 observed meteorology over time. The outermost 10km domain uses initial and boundary conditions from the
614 output of a meteorological simulation, that is itself driven by an analysis field. The outermost domain model then
615 carries out a 36-hour forecast, of which the first 6 hours are discarded as spin-up; the final 30 hours are used as
616 initial and boundary conditions for the rotated 10 km grid cell size domain (the OS10km domain). An OS10km
617 simulation of 30 hours is then carried out. The forecasts run in a repeating cycle from new meteorological analyses
618 on every 36 hours, and hence are constrained by observations to prevent chaotic drift of the forecast over an

619 extended simulation. The outermost 10km domain carries out a 36 hour forecast, of which the first 6 hours is
620 discarded as spin up; the final 30 hours is used as initial and boundary conditions for the rotated 10 km grid cell
621 size domain (the OS10km domain). As noted above, the OS10km domain makes use of a microphysics package
622 matching that of the subsequent higher resolution simulations for better matching of cloud fields at those
623 resolutions. An OS10km simulation of 30 hours is then carried out, with the first 6 hours being discarded as spin-
624 up, and the latter 24 hours forming the initial and boundary conditions for the 2.5 km grid cell size OS2.5km
625 simulation. The OS2.5km simulation is of 24 hours duration. The OS1km simulation covers the same 24 hours (and
626 hence both 2.5km and 1km simulations start from the same OS10km initial conditions at for every 24 hour
627 forecast), with the 2.5km simulation providing boundary conditions thereafter to the OS1km model. Continuity
628 between 24 hour forecasts is thus maintained at the level of the outermost nest. The outermost domain is cycled
629 every 12 hours starting at 0UT and 12UT; however, we have selected the set of contiguous OS2.5km and OS1km
630 24 hour simulations starting from the 12UT continental domain for our comparison.

631 Meteorological boundary conditions for lowest resolution GEM-MACH simulations are taken from operational
632 GEM forecasts, in turn driven by data assimilation analyses performed at the Canadian Meteorological Centre.



633
634 Figure 2. (a) The four nested domains of the GEM-MACH simulations. From outermost to innermost domains,
635 these are CONT10km (outermost, red dots), OS10km (yellow), OS2.5km (green), and OS1km (blue). The model
636 simulations from the two innermost domains are the focus of the present study. (b) Topography in the OS1km
637 domain centred on Fort McMurray, Alberta (m agl). The coloured area corresponds to the central blue domain in
638 (a).

639 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*

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

641 2.3 Model Emissions

642 All emissions data used in this work are described in Zhang *et al.* (2018). These emissions data include (a) direct
 643 observations of stack-specific hourly emissions measured by Continuous Emission Monitoring Systems (CEMS), (b)
 644 regional emissions inventory data from the Cumulative Environmental Management Association (CEMA) - which
 645 had the most detailed stack and process level emission data for the AOSR facilities, including emissions from mine
 646 faces, tailings ponds, and the off-road mining fleet), (c) the 2010 Canadian Air Pollutant Emissions Inventory (APEI)
 647 - which is the most comprehensive national emissions inventory, and which has the largest spatial coverage for
 648 area sources ~~for areas~~ outside the AOSR, and (d) the 2013 National Pollutant Release Inventory (NPRI) (a subset of
 649 the APEI) that is based on emissions reports from large industrial facilities.

650 These emissions data sets primarily describe emissions of pollutants known as criteria-air-contaminants (NO_x ,
 651 VOCs , SO_2 , NH_3 , CO , $\text{PM}_{2.5}$, and PM_{10}) for *major-point sources* (*i.e.*, large emission stacks) and *area sources*. Area
 652 emissions sources typically consist of multiple small mobile sources spread over a large area (*e.g.*, off-road
 653 vehicles), large flux sources such as mine tailings settling ponds or mine faces, and/or large numbers of small
 654 stacks for which no stack characteristic data (volume flow rates, temperatures of emissions, stack diameters),
 655 needed to estimate plume-rise heights, are available.

656 Major-point sources are represented by a single geographical (latitude, longitude) pair of coordinates, and are
 657 assigned to the grid cell in which the point is located. These sources are likely to be the most impacted by model
 658 horizontal grid cell size, as even a large major-point source plume, which in reality may only occupy an emissions
 659 horizontal area on the order of 100 m^2 , is represented by a flux spread over an entire grid cell. A plume from a
 660 major point source within a 2.5km grid cell will thus be immediately diluted to a size of 6.25km^2 upon emission,
 661 whereas the same source with a 1km grid cell will have a cross-sectional horizontal extent of 1km^2 . At the same
 662 time, higher resolution may require a much more accurate representation of model winds close to the sources to
 663 maintain accuracy in evaluation metrics dependant on plume position such as correlation – a wider plume being
 664 more likely to at least partially intersect a monitoring station location than a narrower plume.

665 Area sources that are large compared to both model grid cell sizes (2.5km and 1km) can be expected to be
666 approximated by model grid cells of both resolutions, and are thus expected to be less impacted by model
667 resolution than emissions from point sources. However, smaller area sources (*i.e.* areas intermediate between
668 2.5km and 1km to the side) may be better resolved, and hence have less dilution and higher downwind
669 concentrations, when higher spatial resolution is employed.

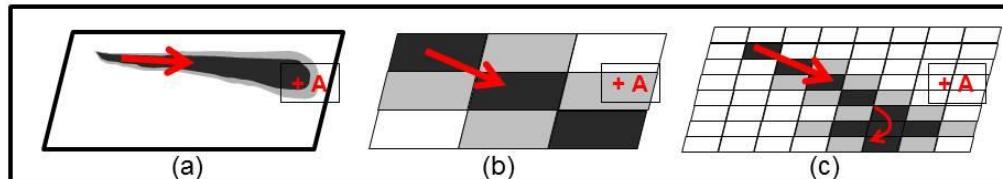
670 In the AOSR, approximately 95% of the SO₂ emissions originate in major-point sources, while NO₂ is
671 apportioned ~40% to major-point sources and ~60% to area sources (Zhang *et al.*, 2018). Consequently our *a*
672 *priori* expectation is that the impact of the resolution change will be strongest for species like SO₂, and less strong
673 for species like NO₂ that are emitted in part by point sources, but may also be apparent for other species and
674 secondary products, such as O₃.

675 [2.41.4](#) Model Evaluation Methodology and Metrics

676 Comparisons between air-quality models and observations usually take the approach of comparing observation
677 and model-generated values paired in time and space, from the observation location and corresponding model
678 grid-cell respectively. We refer to this approach hereafter as our “standard” evaluation, for both 2.5km and 1km
679 simulations. However, we note additional factors aside from grid-cell size may influence the outcome of air-
680 quality model evaluations. For example, the relative skill of the meteorological component of the air-quality
681 model will depend in part on the density of meteorological observation data, incorporated into the model via data
682 assimilation, for the construction of the model’s initial meteorological state. This in turn will influence the local
683 skill of the model’s predicted wind directions and hence the skill of its plume transport. The simulations carried
684 out here focus on the Fort McMurray area, where the nearest available upper air meteorological sounding site is
685 located at the ECCC Stony Plain station, located approximately 500km south-west of the study area. The
686 advantage of higher resolution simulations (*e.g.*, reduced numerical error associated with the discretization of
687 transport operators, and better treatment of local topographic influences) may thus be offset by errors in the
688 predicted *large scale* flow.

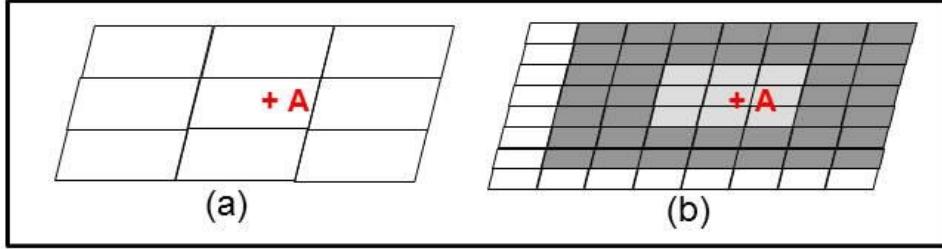
689 While meteorological model synoptic-scale forecast errors may manifest themselves locally as errors in the
690 direction of winds driving local plume transport, other advantages may result from the use of higher resolution
691 air-quality models. Since lower resolution models *de facto* instantaneously redistribute plumes emitted from
692 large stack sources over a larger area, such artificial diffusion will reduce the model’s ability to accurately simulate
693 concentration maxima, and the resulting chemistry, within simulated model plumes. However, the spatial extent
694 of a plume in a model employing a large horizontal grid cell size may be such that its existence may be captured at
695 discrete observing sites. In contrast, forecast plumes in models with smaller horizontal grid cell sizes may
696 correctly capture plume magnitude and chemical behaviour, but may be more subject to errors in the larger scale
697 wind direction. To illustrate this point, Figure 3 shows a conceptual diagram of an actual plume, a large grid cell
698 size model plume, and a small grid cell size model plume, where the latter two simulated plumes are both subject

699 to the same synoptic-scale error in wind forecast direction (indicated by large red arrows; the smaller red arrow in
700 Figure 3(c) indicates the impact of local forcing predicted for the second model). Observation station “+A” is
701 located downwind, and records the presence of the actual plume (Figure 3(a)). The coarse grid cell size simulated
702 plume (Figure 3(b)), despite the error in the forecast wind direction, captures part of the observed plume in the
703 resulting time series at the observation station location. In contrast, the small grid cell size plume (Figure 3(c)),
704 despite resolving the plume shape (and plume-internal chemistry) to a greater degree than the coarse grid cell size
705 simulated plume, fails to record the presence of the plume at the observation location. A simple paired
706 observation-model time series evaluation would thus suggest that the former model has superior performance to
707 the latter model in this example, despite the latter model having created a more “realistic” plume in terms of the
708 maximum concentration reached, albeit in the wrong location, due to synoptic-scale forecast wind direction error.
709 In this particular instance, the magnitude of the smaller grid cell size simulated plume is more realistic than that of
710 the coarse grid cell size plume, but this improvement will not be captured in a standard evaluation analysis. Shifts
711 in plume location across individual grid cells away from the location of an *in-situ* observation are more likely grid
712 cell size decreases. In this example, a standard analysis would impose a more stringent expectation on the smaller
713 grid cell size simulation to correctly identify plume locations.



714
715 Figure 3. Schematic comparison of surface concentration contours and model grid cell values of a transported pollutant
716 plume from a large stack (termed a “point” source). Wind direction shown by red arrows. Monitoring station location
717 marked by “+A”. (a) Actual plume. (b) Coarse grid cell size air-quality model prediction. (c) Fine grid cell size air-quality model
718 prediction. Note the change in wind direction between observations (a) and simulations (b,c) associated with errors in the
719 forecast of the synoptic wind.

720 In addition to the standard analysis, we perform additional analyses that examine the model’s ability to resolve
721 plumes in the vicinity of the observation station, in order to attempt to evaluate the potential for higher
722 resolution simulations to provide potential benefits which that are may be masked by synoptic scale forcing
723 errors, in addition to the standard analysis, we perform additional analyses that examine the model’s ability to
724 resolve plumes in the vicinity of the observation station. This strategy is illustrated in Figure 4.



725

726 Figure 4. Scale diagram of the same region in (a) 2.5km grid cell size simulation and (b) a 1km grid cell size simulation.
 727 Region enclosed by light grey / dark grey shading in (b) represents the nearest nine / forty-nine 1km gridpoints surrounding
 728 the observation location "A".

729 Figure 4(a) shows an observation station enclosing the nine nearest-neighbour model grid-cells for a 2.5km grid
 730 cell size, while Figure 4(b) shows the corresponding 1 km grid cell size map, with the nine nearest-neighbour
 731 model grid-cells shown in light grey, the forty-nine nearest grid cells shown in the region enclosed in dark grey.
 732 Figure 4(a) encloses a region of 56.25 km^2 ($7.5 \times 7.5 \text{ km}$), while the light grey region in Figure 4(b) encloses 9km^2 ,
 733 and the darker grey region encloses 49 km^2 .

734 As noted above, in a formal mathematical sense, the smallest region resolvable by an Eulerian grid model is twice
 735 the size of the model grid cell size (relating to the Nyquist frequency of the model); hence the smallest resolvable
 736 feature spans two model grid cells in each direction. However, in a practical sense, a total of nine grid cells
 737 centred on the observation station must be used to allow a boundary of two grid cells in any direction. Sampling
 738 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
 739 events occurring at discrete location "+A". The closest corresponding sampling region available to the 1 km model
 740 (Figure 4(b)) is shown in dark grey. The light grey region of Figure 4(b) represents the closest 1 km grid cell size
 741 region that corresponds to the single 2.5 km grid cell in which the observation station is located in Figure 4(a). We
 742 attempt to ascertain model performance in these approximately equivalent regions around each observation
 743 station, in the analysis that follows.

744 Our approach follows two steps:

745 (1) From the 2.5km simulation, in addition to the predicted model value at the grid-cell containing the
 746 observation location, we determine the model grid-cell value in the nine grid-cells surrounding the
 747 observation station location which has the closest value to that observed at the station. This represents the
 748 model's "best estimate" of the value at the observation station location itself, to the model's ability to resolve
 749 features at 2.5km grid cell size.

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

751 the observation value from: (a) the nearest nine grid-cells to the observation station location, and (b) the
752 nearest 49 grid-cells to the observation station location. The former represents the model's "best estimate"
753 of the value at the observation station location itself, while the latter represents the 1km model's best
754 estimate in the closest equivalent region to the limiting resolution of the 2.5km model.

755 Comparing the resulting statistical measures of each of these selected values with observations, in addition to the
756 standard analysis, thus evaluates the model's best attempt to resolve features for the specified grid cell size, and
757 allows cross-comparison of model performance within nearly equivalent areas. Cross-comparing the statistical
758 values for the different regions described above shows the model's ability to resolve features such as plumes from
759 the standpoint of the region represented at the different grid cell sizes. If synoptic-scale transport direction errors
760 creates situations similar to that depicted in Figure 3(a), a standard comparison of error would be expected to
761 show little benefit to higher resolution. However, the "best model estimate" comparisons would capture the
762 ability of the higher resolution model to more accurately simulate the magnitude of the plume, if not its spatial
763 location. Each of these selection procedures will be employed in the surface concentration comparisons which
764 follow.

765 We evaluate our model simulations against observations made at surface monitoring networks in the vicinity of
766 the Athabasca oil sands, and aboard an instrumented aircraft, the National Research Council of Canada Convair.
767 For the surface monitoring data, hourly time series of model output were matched to station time series using the
768 different strategies described above. For the aircraft observations, we extract model values through temporal and
769 spatial interpolation to the aircraft's position during the flights and only perform the standard analysis, as well as
770 examining the behaviour of the two simulations along cross-sections corresponding to the flight paths.

771 Our statistical metrics for evaluation are common to many other air-quality applications, and were computed
772 using the 'modstat' function from the OpenAir R package (Carslaw and Ropkins, 2012). [Further discussion of](#)
773 [different metrics for model evaluation may also be found in Yu et al., \(2006\).](#) The statistics calculated here
774 include: mean bias (MB; perfect score: zero), mean absolute gross error (MGE; perfect score: zero), normalised
775 mean bias (NMB; perfect score: zero), normalised mean gross error (NMGE: perfect score: zero), root mean
776 squared error (RMSE; perfect score: zero), correlation coefficient (r , perfect score: unity), coefficient of
777 efficiency (COE: a perfect score is unity, a zero/negative score means the model is equivalent/less predictive
778 than the mean of the observations), and the index of agreement (IoA; perfect agreement is unity, and -1
779 indicates no agreement or little variability).

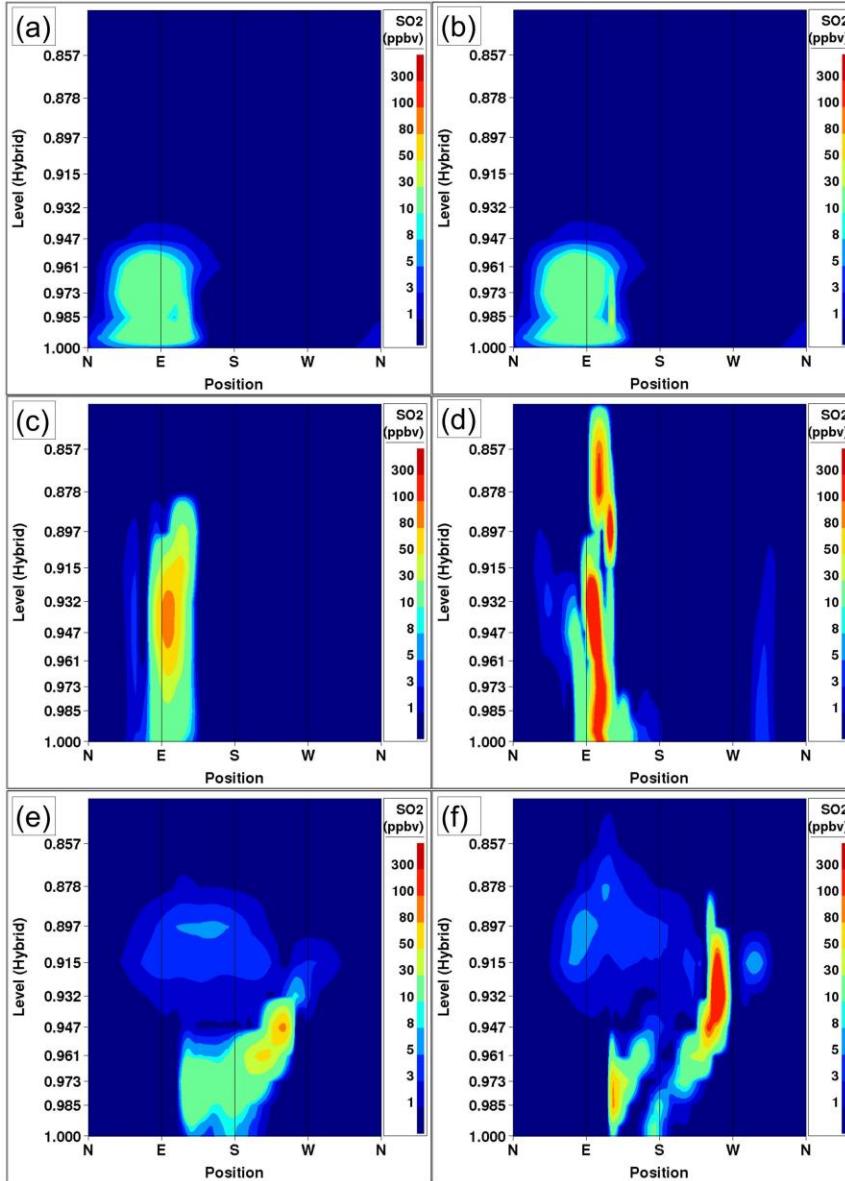
780 [3.2](#) Simulation Comparisons and Evaluation

781 3.1 Model-to-model comparisons and averages

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

784 domain, in order to provide a qualitative comparison of the differences in simulations for the two simulations, and
785 then continue with the quantitative comparison. Figure 5 compares OS2.5km (left column) and OS1km (right
786 column) simulation results for a cross-section located 0.2km from a major SO₂ emissions source at 0, 12 and 24
787 hours into a given simulation day.

788 The model results are identical at hour 0 due to both the OS2.5km and OS1km models being initialized from the
789 OS10km data at this time (small differences in Figure 5(a,b) are due to slight mis-matches in the cross-section
790 locations). Subsequent cross-sections show the OS1km model is capable of resolving both higher absolute mixing
791 ratio values, and sharper gradients, within 12 hours of simulation time (Figure 5 (c,d)). Multiple plumes are
792 resolved by 12 hours of simulation time in the 1km grid cell size simulation, along with markedly different plume
793 heights, plume structure, and a factor of two increase in the magnitude of plume mixing ratios relative to the
794 lower grid cell size simulation, and these differences persist into the 24th simulation hour (Figure 5(e,f)). Mixing
795 ratio differences of these magnitudes are to be expected given the increase in resolution, but Figure 5 shows that
796 other important aspects of the predicted plumes have changed. The plume heights are a function of predicted
797 local stability conditions in the grid-square containing the source, and the variation shown here represents a
798 substantial change in the predicted local stability for the origin sources of these plumes, resulting from the change
799 in model horizontal grid cell size.



800

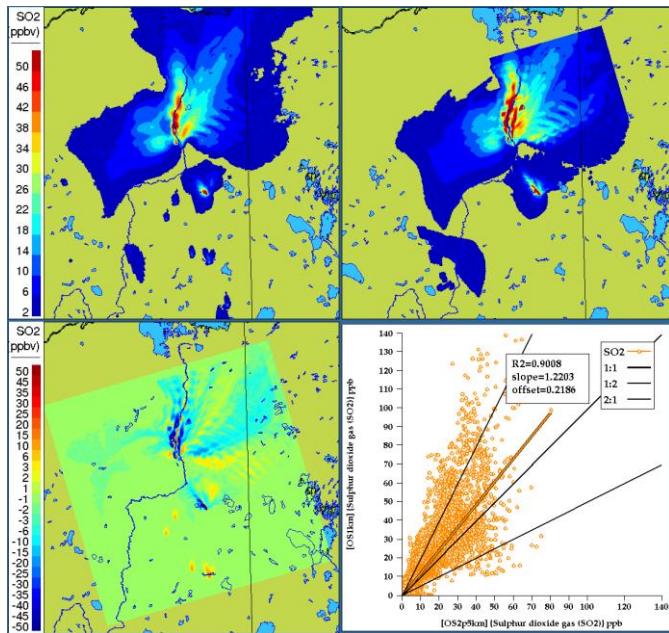
801 Figure 5. Comparison of simulated SO₂ plume mixing ratios (ppbv) located 0.2km from a major point source, for OS2.5km
 802 simulations (left column) and OS1km simulations (right column), at 0 (a,b), 12 (c,d), and 24 (e,f) hours into a 24 hour
 803 simulation.

804

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

805 OS10km data at this time (small differences in Figure 5(a,b) are due to slight mis-matches in the cross-section
 806 locations). Subsequent cross sections show the OS1km model is capable of resolving both higher absolute mixing
 807 ratio values, and sharper gradients, within 12 hours of simulation time (Figure 5 (c,d)). Multiple plumes are
 808 resolved by 12 hours of simulation time in the 1km grid cell size simulation, along with markedly different plume
 809 heights, plume structure, and a factor of two increase in the magnitude of plume mixing ratios relative to the
 810 lower grid cell size simulation, and these differences persist into the 24th simulation hour (Figure 5(e,f)). Mixing
 811 ratio differences of these magnitudes are to be expected given the increase in resolution, but Figure 5 shows that
 812 other important aspects of the predicted plumes have changed. The plume heights are a function of predicted
 813 local stability conditions in the grid square containing the source, and the variation shown here represents a
 814 substantial change in the predicted local stability for the origin sources of these plumes, resulting from the change
 815 in model horizontal grid cell size.

816 Figure 6 compares the maximum surface SO₂ during the entire period for each simulation, as well as the difference
 817 in maximum SO₂ between the simulations, along with a scatterplot of OS2.5km versus OS1km simulation results.
 818 (where, in the latter two panels, OS2.5km values were assigned to the corresponding OS1km grid-cell locations
 819 using the nearest-neighbour approach).



820
 821 Figure 6. Comparison of total-simulation *maximum* surface SO₂ mixing ratios (ppbv) at (a) 2.5km and (b) 1km grid cell size
 822 (ppbv). (c) Difference (2.5km – 1km). (d) Scatterplot of 2.5km (x-axis) versus 1km (y-axis) total simulation average grid-cell

823 surface SO₂ mixing ratios.

824 The maximum surface concentrations tend to show more elongated structures at the smaller grid cell size, ^{2.5}
825 comparing Figures 6(a,b), particularly for plumes in the western (left) half of the OS1km domain. The difference
826 plot (Figure 6(c)) shows that local maximum concentration differences of up to -45 ppbv occur, due to changes in
827 the placement and maximum concentration of high concentration plumes. The scatterplot of Figure 6(d) shows
828 that OS1km model has a demonstrated ability to achieve higher concentrations than the OS2.5km model, with a
829 slope of 1.22, and a noticeable clustering of values along the 1:2 line. While these results are not unexpected
830 since approximately 95% of the SO₂ emissions in the domain originate in large stack, or point, sources, and hence
831 initial concentrations at source would be expected to 6.25x higher in the OS1km simulation, they also suggest that
832 a substantial improvement in the OS1km model's ability to capture SO₂ concentrations *should* be possible. That is,
833 the results of the two models are substantially different, and given the reduction in numerical error expected with
834 employing a smaller grid cell size, the latter might be expected to outperform a larger grid cell size model.
835 However, as we shall demonstrate in the next section, plume placement errors such as depicted in Figure 3 play a
836 substantial role in model performance as grid cell size decreases.

837 3.2 Quantitative comparisons

838 3.2.1 Surface observation comparison

840 The locations of the local network of 10 surface monitoring stations located near the sources of emissions in the
841 region (oil sands facilities) are shown in Figure 7. As noted in section 2.4, we carry out several analyses:

842 (1) The standard evaluation (model values are extracted from the model grid-cells containing the observation
843 stations, at both grid cell sizes).

844 (2) Equal areas of representativeness, 1km and 2.5km grid cell sizes (the nearest nine OS1km grid cells are
845 compared to the OS2.5km single cell evaluation in two ways):

846 a. Averaging of the OS1km results across the nine grid cells prior to evaluation (to determine whether
847 the mean value is better represented by the smaller grid cell size, similar to the approach taken in
848 Kang *et al.* (2007)).

849 b. Selection of the *best* of the nine grid cells (closest to the observation value), to determine the extent
850 to which the OS1km model is capable of better representing the concentrations somewhere within
851 the corresponding OS2.5km model grid cell, if not at the OS1km cell closest to the observation
852 location. Higher scores for the 1km grid cell size simulation in this case would indicate that while
853 errors in plume positioning (for example due to errors in the synoptic scale flow) negate some of the
854 advantages of the OS1km simulation, the plume may be better represented by the OS1km simulation
855 within the 2.5km grid cell's area.

856 (3) Equal areas of representativeness and equal regions of variability (nearest nine 2.5km cells are compared to
857 the nearest forty-nine 1km cells). Here we make the assumption that the 2.5km grid cell size model's ability
858 to resolve features is limited to the surrounding three grid cells in each horizontal dimension, and make use of
859 the closest-in-size block of corresponding 1km cells (a 7×7 grid centered on the cell containing the
860 observation point). In both cases, the model value closest to the observations is chosen prior to evaluation.

861 While evaluations (2b) and (3) deliberately select the “best” value, they also provide a quantitative estimate of
862 the extent to which each model is capable of achieving the correct answer within roughly equal representative
863 areas centered on the observation station locations. These comparisons are intended to evaluate (a) the
864 extent to which the 1km grid cell size is capable of improving simulation results despite, *e.g.*, the larger scale
865 flow resulting in errors in the plume placement, and (b) whether the 1km grid cell size model is capable of
866 outperforming the 2.5km grid cell size model *over equivalent regions*. In the last test, we place both models on
867 an equal footing with regards to the region being represented, as well with regards to allowing cell-to-cell
868 variability and the selection of a closest match to observations.

869 Our evaluation is presented as tables of statistical metrics. The comparisons employing the nearest neighbour
870 approach are described with a “B#” superscript suffix, denoting that the “best” sample within a square centred
871 on the observation point containing a total of # grid cells (*e.g.* the OS1km⁸⁹ label denotes a comparison
872 between observed data and the simulation grid cell within a 3×3 grid-cell square centered about the
873 observation point). Similarly, an A# superscript describes a comparison between the observations and the
874 Average of the # square of grid cells centered on the observation point.

875 Comparisons to surface concentrations were performed using publicly available data collected by the Wood
876 Buffalo Environmental Association (WBEA), which operates the air-quality monitoring network residing within
877 the OS1km domain. The monitoring station locations are shown in Figure 7. The statistical performance of the
878 models, calculated using the procedure outlined above, are given in Tables 2 through 5, for SO₂, NO_x, O₃, and
879 PM_{2.5}, respectively.



880

881 Figure 7. Illustration of the OS1km domain, with observation measurement comparison for SO₂, and NO_x (first two columns,
 882 view of station locations. Monitoring stations are shown as purple outline squares in both images. Light grey
 883 regions in the background satellite image (b) are oil sands open-pit mining operations.

884 In the *standard* model grid cell to observation measurement comparison for SO₂, and NO_x (first two columns,
 885 Tables 2 and 3), the OS1km simulation had *worse* scores for all the metrics considered here. For O₃, the OS1km
 886 model had the better score for the correlation coefficient and root mean square error, and worse scores for all
 887 remaining model evaluation metrics. For PM_{2.5}, the OS1km model had higher performance for the correlation
 888 coefficient and biases, while the OS2.5km model outperforms the OS1km model for all other metrics examined
 889 here. Based on a standard analysis, the OS1km model thus performs poorly compared to the OS2.5km model; the
 890 expected advantages associated with reduced numerical error in transport at smaller grid cell sizes are being offset
 891 by other factors controlling the net model error.

892 When *the* standard evaluation is compared to the *average* of the nearest nine 1km simulation grid cells
 893 surrounding the observation point (*first three columnsthird column* of the tables), an intermediate result appears.
 894 For SO₂ (Table 2) the nine-cell OS1km average has the best performance for correlation coefficient - indicating a
 895 better time distribution of events may be achieved by a nine cell average at 1km grid cell size. The other metrics for
 896 the A9 simulation are intermediate between the two standard evaluations for each simulation, indicating that some
 897 of the performance loss resulting from the use of 1km grid cell size is reduced through averaging results to
 898 approximately the same size regions as the OS2.5km grid cell size. The latter result holds for all metrics for NO_x
 899 (including R, see Table 3). For ozone (Table 4), averaging the nine nearest OS1km grid cells prior to measurement
 900 gives the best performance for R and RMSE, and worse performance for the other metrics. For PM_{2.5} (Table 5), all
 901 metrics for the OS1km nine grid-cell average aside from the bias fall mid-way between the two standard
 902 methodology evaluations. Averaging the smaller grid cell size model results thus shows a marginal improvement,

903 depending on the species, but overall does not compensate for the decrease in performance resulting from going
904 to the smaller grid cell size.

905 We next ask the question, “Does a more accurate simulation value *exist* within the same region of the 1km model
906 as is encompassed by a 2.5km grid cell?” (fourth column of these Tables), by selecting the model value in the
907 nearest nine 1km grid cells with the closest match to observations and comparing to the corresponding single 2.5
908 km_grid cell. A dramatic improvement in the relative OS1km performance metric scores can be seen. For each of
909 Tables 2 through 5, this “best of nine” 1km comparison outperforms the previous 3 comparisons (columns 1
910 through 3), for all metrics. These improvements are sometimes dramatic (e.g. a doubling of correlation coefficient
911 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
912 of coefficient of error from negative to positive values for O₃, and a reduction in the coefficient of error for PM_{2.5} by
913 a factor of 2.5 compared to the nearest competing value from the previous evaluations. The coefficient of
914 efficiency for SO₂ and O₃ make the transition from negative to positive values when the “best-of-nine” methodology
915 is used, indicating that the model is able to better predict the observations than the observed mean, somewhere
916 within an equivalent area. This evaluation suggests that the OS1km model does *contain* a better result within the
917 same approximate region encompassed by a 2.5km grid cell. However, the location of that better result may be
918 subject to positioning error, such as described in Figure 3.

919 A valid argument could be made that the methodology employed in this fourth evaluation is subject to selection
920 bias, in that the selection of a *best* value in the case of the nearest nine 1km simulation places that model
921 simulation at an advantage relative to the 2.5km model. To address this last issue, the final two additional
922 methodologies for evaluation were employed, still maintaining the same approximate area of representativeness
923 for a grid cell, namely choosing the best value out of the nearest *nine* 2.5km grid cells (the limiting resolution of this
924 model simulation), and the best value out of the nearest *forty-nine* 1km grid cells (fifth and sixth columns of Tables
925 2 through 5, respectively). That is, we attempt to place the two models on an equal basis with regards to selection
926 bias within a given region containing an observation station.

927 Two important results can be seen from this final evaluation. First, as was the case for the “Best of 9” for the
928 OS1km simulation compared to the standard OS1km evaluation, the “Best of 9” for the OS2.5km simulation has a
929 considerably better performance than the standard OS2.5km evaluation (compare fifth and first columns, Tables 2
930 through 5). That is, the OS2.5km model may *also* be subject to location errors in transported species representation
931 which influence model performance. However, when performance within the 56.25 km² area surrounding each
932 measurement point in the OS2.5km “Best of 9” evaluation is compared to the 49 km² area surrounding the
933 measurement points in the OS1km “Best of 49” simulation (*i.e.* compare columns five and six in Tables 2 through 5),
934 it can be seen that the OS1km model outperforms the OS2.5km model for all metrics for O₃, and PM_{2.5}, and all
935 metrics aside from bias for SO₂ and NO_x. That is, despite the OS1km model having a slight disadvantage in the

936 relative size of the representative area containing the measurement station location, and both models being
937 allowed a similar selection strategy, the OS1km model is capable of generating values closer to the observations
938 than the OS2.5km model within an equivalent sub-region, across most of the metrics and chemical species
939 considered here.

940 This final result is strongly suggestive of the presence of issues such as illustrated in Figure 3. These may include
941 errors in the larger scale synoptic wind flow, combined with the reduced size of plumes as grid cell size is reduced,
942 leading to more “misses” than “hits” for a given recorded event at a measurement station compared to the coarse
943 grid cell size model. There may be multiple additional causes for such errors (examples include poor observation
944 density in the region for model initialization, underlying lower resolution boundary condition fields such as
945 topography not improving with the reduction in grid cell size, inaccuracies in land use fields used in meteorological
946 modelling due to rapid development, and errors in other aspects of the reaction transport modelling system aside
947 from horizontal resolution). The expected advantages of the small grid cell size, such as better representation of
948 the concentrations of species within plumes and hence better representation of their reactive chemistry (c.f.
949 Lonsdale *et al.*, 2012), may be lost in a standard performance analysis due to these other issues.

950 Our analysis suggests that a practical limit in the benefits of increasing model accuracy may be reached when
951 resolution exceeds some threshold, as a result of other errors inherent in the modelling system. However, the
952 analysis also suggests that if these non-resolution-related errors are corrected, the benefits of adopting a smaller
953 grid cell size may be substantial. For example, meteorological data assimilation employing a dense monitoring
954 network for a specific area of interest would be expected to show a greater impact for smaller than larger grid cell
955 sizes, due to the greater ability of the former to take advantage of the observation density in correcting the initial
956 meteorological state. We note that recent work applying land use data assimilation (Carrera *et al.*, 2015) to
957 regional 2.5km grid cell size weather simulations (Milbrandt *et al.*, 2016) have suggested that such data assimilation
958 may indeed improve forecast skill at the very local scale.

959 Table 2. Surface SO₂ observations to model comparison for entire simulation period (ppbv)

Evaluation Metric	OS2.5km	OS1km	OS1km ^{A9}	OS1km ^{B9}	OS2.5km ^{B9}	OS1km ^{B49}
Index of Agreement	0.237	0.154	0.207	0.601	0.701	0.810
Pearson Correlation Coefficient	0.290	0.230	0.295	0.604	0.672	0.848
Normalized Mean Gross Error	2.128	2.363	2.212	1.114	0.834	0.529
Mean Gross Error	2.918	3.240	3.034	1.528	1.143	0.725

Formatted: Font: Bold

Formatted Table

<u>CeE</u> Coefficient of Error	-0.525	-0.693	-0.585	0.202	0.403	0.621
<u>RMSE</u> Root Mean Square Error	7.063	9.665	7.876	4.436	3.671	2.618
<u>NMB</u> Normalized Mean Bias	1.130	1.376	1.299	0.347	-0.010	0.017
<u>MB</u> Mean Bias	1.550	1.887	1.781	0.475	-0.013	0.024

- 5466 Samples used

Formatted: Justified, Bulleted + Level: 1 + Aligned at: 9.52 cm + Indent at: 10.16 cm

960

961 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}
<u>Index of Agreement</u> ^{10A}	0.177	0.138	0.152	0.416	0.589	0.665
<u>Pearson Correlation Coefficient</u> ^{10B}	0.143	0.114	0.116	0.165	0.305	0.388
<u>Normalized Mean Gross Error</u> ^{NGME}	1.520	1.593	1.567	1.079	0.760	0.619
<u>Mean Gross Error</u> ^{GME}	12.898	13.518	13.296	9.156	6.447	5.255
<u>Coefficient of Error</u> ^{CeE}	-0.646	-0.725	-0.697	-0.168	0.177	0.329
<u>Root Mean Square Error</u> ^{RMSE}	28.052	35.197	34.644	25.782	15.315	13.704
<u>Normalized Mean Bias</u> ^{NMB}	0.493	0.570	0.542	0.174	-0.027	-0.063
<u>Mean Bias</u> ^{MB}	4.183	4.834	4.597	1.477	-0.231	-0.531

- 3257 Samples used

962

963 Table 4. Surface O₃ observations to model comparison for entire simulation period (ppbv)

Evaluation Metric	OS2.5km	OS1km	OS1km ^{A9}	OS1km ^{B9}	OS2.5km ^{B9}	OS1km ^{B49}
<u>Index of Agreement</u>	0.414	0.405	0.404	0.527	0.637	0.690
<u>Pearson Correlation Coefficient</u>	0.496	0.506	0.515	0.606	0.688	0.738
<u>Normalized Mean Gross Error</u>	0.660	0.670	0.672	0.534	0.410	0.349
<u>Mean Gross Error</u>	10.757	10.915	10.949	8.692	6.673	5.687
<u>Coefficient of Error</u>	-0.172	-0.189	-0.193	0.053	0.273	0.380

<u>Root Mean Square Error</u>	<u>16.040</u>	<u>15.859</u>	<u>15.794</u>	<u>13.305</u>	<u>11.084</u>	<u>9.719</u>
<u>Normalized Mean Bias</u>	<u>0.527</u>	<u>0.559</u>	<u>0.579</u>	<u>0.463</u>	<u>0.337</u>	<u>0.304</u>
<u>Mean Bias</u>	<u>8.579</u>	<u>9.104</u>	<u>9.431</u>	<u>7.536</u>	<u>5.488</u>	<u>4.945</u>

• 2189 Samples used

Table 4. Surface O₃ 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
f	0.496	0.506	0.515	0.606	0.688	0.728
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
CeE	-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

• 2189 Samples used

Table 5. Surface PM_{2.5} observations to model comparison for entire simulation period (μg m⁻³)

Evaluation Metric	OS2.5km	OS1km	OS1km ^{A9}	OS1km ^{B9}	OS2.5km ^{B9}	OS1km ^{B49}
<u>Index of Agreement</u> IoA	0.280	0.262	0.267	0.412	0.508	0.572
<u>Pearson Correlation Coefficient</u> f	0.201	0.216	0.214	0.314	0.376	0.466
<u>Normalized Mean Gross Error</u> NGME	0.791	0.811	0.806	0.647	0.541	0.471
<u>Mean Gross Error</u> GME	5.342	5.478	5.441	4.365	3.651	3.181
<u>Coefficient of Error</u> CeE	-0.439	-0.476	-0.466	-0.176	0.016	0.143

<u>Root Mean Square Error</u> <u>RMSE</u>	8.286	8.786	8.663	7.117	6.169	5.690
<u>Normalized Mean Bias</u> <u>NMB</u>	-0.268	-0.257	-0.257	-0.289	-0.299	-0.287
<u>Mean Bias</u> <u>MB</u>	-1.812	-1.734	-1.736	-1.948	-2.016	-1.937

969

- 3377 Samples used

970 The surface observation data were also analyzed by time-of-day, with both observations and simulations split into
 971 daytime (hours 9:00 to 18:00 local time) and nighttime (hour 19:00 to 8:00 local time) data pairs (Appendix, Tables
 972 A1 through A8, [Carslaw and Ropkins, 2012](#)). Within each of these diurnally segregated time periods, the broad
 973 aspects of the comparison were the same as for the “all data” Tables 2 to 5 above: the OS1km simulations tended
 974 to have reduced performance in a standard analysis, averaging improved but not completely ameliorated the
 975 performance of the OS1km simulation, a methodology employing the best of nine OS1km grid cells had superior
 976 performance to the two standard comparisons, and comparison of the “best of” methodologies for equal areas
 977 showed better performance for the OS1km compared to the OS2.5km simulation. We also noted substantial
 978 differences in the day and night performance of both models across the methodologies. For example, daytime SO₂
 979 and NO_x performance within a given model and comparison methodology was usually better than nighttime
 980 performance for IOA,R, NGMEMGE, COE and NMB, while worse for RMSE, while nighttime O₃ performance was
 981 better for IOA, r, NGMEMGE, and COE. Daytime PM_{2.5} performance was better than nighttime for IOA, r, COE, and
 982 NMB. [The study area is located in a broad river valley with frequent slope-defined anabatic/akatabic and drainage](#)
 983 [flow events. These often have a diurnal nature, and may explain part of the day/night differences. Example](#)
 984 [sources of these differences may include the relative ability of the driving meteorological model to capture daytime](#)
 985 [versus nighttime mixed layer turbulence and the planetary boundary layer height.](#)

986 3.2.2 Comparisons to Aircraft Observations

987 Twenty-two aircraft observation flights were carried out during the study simulation period – we present
 988 statistical comparisons using the standard approach only, here (model grid cell containing the observation point to
 989 observation data at the aircraft location). Model values were linearly interpolated in time and space to the
 990 aircraft observation locations and times (aircraft observations were on a 10s interval.) We begin with a composite
 991 comparison across all observation times, in Table 6.

992 Table 6. Aircraft observation comparisons, SO₂ and NO₂ (ppbv)

SO ₂ (21787 samples)		NO ₂ (18310 samples)	
OS2.5km	OS1km	OS2.5km	OS1km

Formatted Table

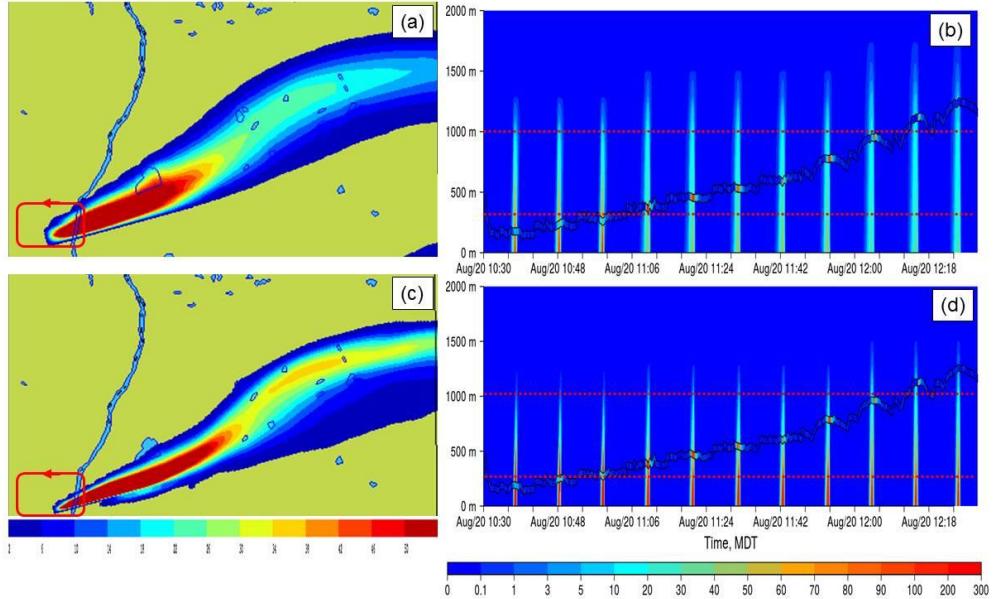
<u>Index of Agreement</u> ^{loA}	0.63	0.62	0.61	0.58
<u>Pearson Correlation Coefficient</u> ^r	0.26	0.28	0.39	0.34
<u>Normalized Mean Gross Error</u> ^{NGME}	1.07	1.09	0.90	0.96
<u>Mean Gross Error</u> ^{GME}	3.98	4.06	1.56	1.68
<u>Coefficient of Error</u> ^{CoE}	0.27	0.25	0.23	0.17
<u>Root Mean Square Error</u> ^{RMSE}	12.84	13.97	3.12	3.62
<u>Normalized Mean Bias</u> ^{NMB}	-0.31	-0.29	-0.26	-0.20
<u>Mean Bias</u> ^{MB}	-1.17	-1.07	-0.45	-0.34

994 The results are in general similar to the surface analysis, in that the OS1km simulation tended to have worse
 995 performance than the OS2.5km simulation (exceptions being the biases for both SO₂ and NO₂, and the slightly
 996 better OS1km correlation coefficient for SO₂). One striking difference between the first two columns of Tables 2
 997 and 3 and Table 14 are the magnitude of the differences between the simulations. Aloft (Table 6), the differences
 998 in performance metric magnitudes between OS2.5km and OS1km simulations are much smaller than at the
 999 surface (Tables 3 and 4). The biases are negative aloft, while positive at the surface, indicating that both models
 1000 may be lofting plumes to insufficient distances; one of the possible (non-horizontal grid cell size dependent)
 1001 causes of model error may be in the extent of vertical transport. This possibility is examined in more detail in
 1002 Akingunola *et al.* (2018, and Gordon *et al.* (2018). An example of this behaviour is shown in Figure 8; both
 1003 plumes fumigate to the surface, while the observed plume resides largely aloft. The OS1km model captures the
 1004 higher concentrations to a better degree, but the impact of excessive fumigation more than offsets this
 1005 improvement, as is shown by the performance evaluation of Table 7, where both models have negative biases
 1006 aloft. In this particular case, the tendency of the model to overestimate the extent of fumigation has a bigger
 1007 impact on performance than grid cell size. Garcia-Menendez *et al.* (2014) have noted similar results for forest fire
 1008 plume prediction.

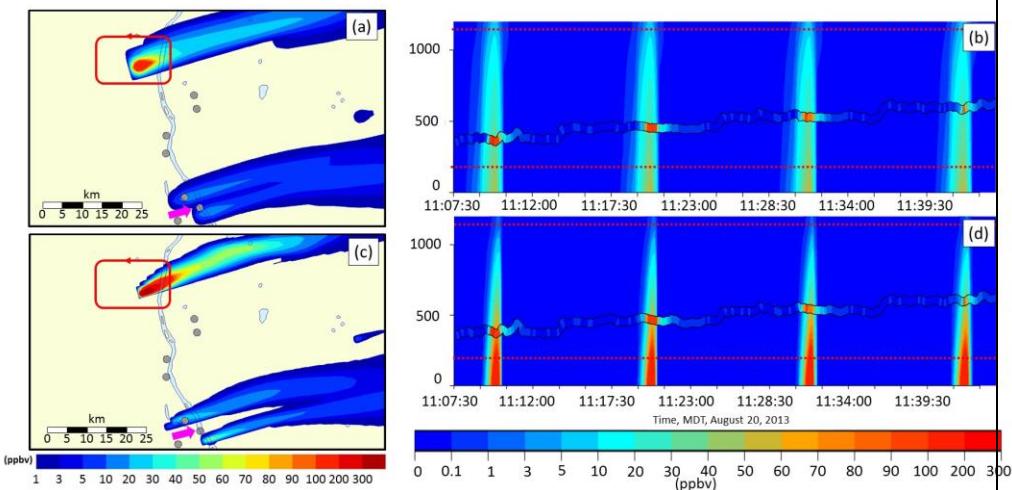
Formatted: Font: Italic

1009 Panels (a) and (c) of Figure 8 provide a further example of the kind of situation referenced in Figure 3; surface
 1010 monitoring station locations are depicted as grey circles, one of which is identified with a pink arrow. This station
 1011 lies within the plume at 2.5km resolution (Figure 8(a)), and outside of the plume at 1km resolution (Figure 8(c)).
 1012 While the plume direction is the same at both scales, that is, the large-scale wind field controls the positioning of
 1013 the plume axis, the smaller grid cell size simulation places a stronger constraint on the accuracy of the wind field.
 1014 For example, if the simulated large-scale flow direction was inaccurately predicted by only a few degrees, the
 1015 plume would not appear in the 1km simulation time series at this location, while registering as present in the
 1016 2.5km simulation. Nevertheless, the plume maximum concentration is better captured by the smaller grid cell size
 1017 simulation (compare maximum values in observed aircraft SO₂, Figure 8 (b, d)). The higher resolution simulation
 1018 may thus more accurately simulate the plume maximum concentration – but not its placement in space, as was
 1019 hypothesized in Figure 3.

1020



1021



1022

1023 Figure 8. Comparison between OS2.5km ([top row a,b](#)) and OS1km ([bottom row c,d](#)) simulations for SO₂ relative to
1024 aircraft observations (ppbv). ([a,c](#)): Simulated surface concentrations of SO₂, with the flight track shown as a red line.
1025 Grey circles: surface monitoring station locations; pink arrow indicates a station located inside a plume at
1026 2.5km resolution ([a](#)), and outside the plume at 1km resolution ([c](#)). ([b,d](#)): Portion of the simulated concentration
1027 profiles along the flight path as a function of time, with the successive intersections of the flight path with the
1028 plume appear as background colour contours. Observed SO₂ aboard the aircraft are shown between the two
1029 black lines. Vertical axis is elevation above the ground, which show the elevation of the aircraft; the aircraft

1030 elevation is increasing with successive passes on successive passes around the facility. Dotted lines show the
1031 upper and lower vertical extent of the observed plume;— Note—note that for both model simulations, the plume
1032 erroneously fumigates the surface._

1034 Table 7. Standard performance evaluation of Flight 8 for SO₂ (ppbv)

	OS2.5km	OS1km
<u>Index of Agreement</u> ^{IoA}	0.69	0.68
<u>Pearson Correlation Coefficient</u> ^r	0.42	0.31
<u>Normalized Mean Gross Error</u> ^{NGME}	1.04	1.09
<u>Mean Gross Error</u> ^{GME}	4.02	4.25
<u>Coefficient of Error</u> ^{CoE}	0.39	0.35
<u>Root Mean Square Error</u> ^{RMSE}	16.72	20.57
<u>Normalized Mean Bias</u> ^{NMB}	-0.42	-0.34
<u>Mean Bias</u> ^{MB}	-1.63	-1.32

1035 Formatted Table

1261 samples used.

1036 Meanwhile other flights show a clear advantage of the OS1km model. One example is given by the NO₂
 1037 performance evaluation of Table 8 and depicted in Figure 9, for Flight 17 (a similar flight plan carried out around
 1038 the same facility as Flight 8). While the correlation coefficient degraded slightly in the OS1km resolution
 1039 simulation, all other performance measures were improved with the decrease in grid cell size. Two time versus
 1040 height profile cross-sections for Flight 17 are shown in Figure 9. In the upper two panels, the OS2.5km (Figure
 1041 9(a)) and OS1km (Figure 9(b)) simulations are compared for the portion of the overall flight track circling the given
 1042 facility. This comparison clearly shows that the OS1km model does a better job of capturing the width of the high
 1043 concentration region of the plume – however, the location of the model plume lags the observations. During this
 1044 portion of the flight alone, the OS2.5km model statistics, particularly the correlation coefficient, outperform the
 1045 OS1km model, due to this issue of plume location mismatching. Figures 9(a,b) may be compared to Figure 3(a,b) –
 1046 the same situation is depicted in both Figures. Figure 9(c,d) show the OS2.5km simulation (10(c)) and OS1km
 1047 simulation results in another portion of the flight – here the OS1km performance for most statistics was better
 1048 than the OS2.5km model performance. The OS1km model (Figure 9(d)) captures the existence of a lower
 1049 concentration layer aloft in the right-hand side of the cross-section, and the existence of low concentration
 1050 intervening layers, as well as the overall lower concentrations of SO₂, while the OS2.5km model does not resolve
 1051 these fine scale and lower concentration features. We note here that IoA, CoE and the other error measures
 1052 capture the visual impression that the OS1km model outperforms the OS2.5km model for this flight, while the
 1053 correlation coefficient is highly dependent on the placement of the plume maximum in the upper two panels.

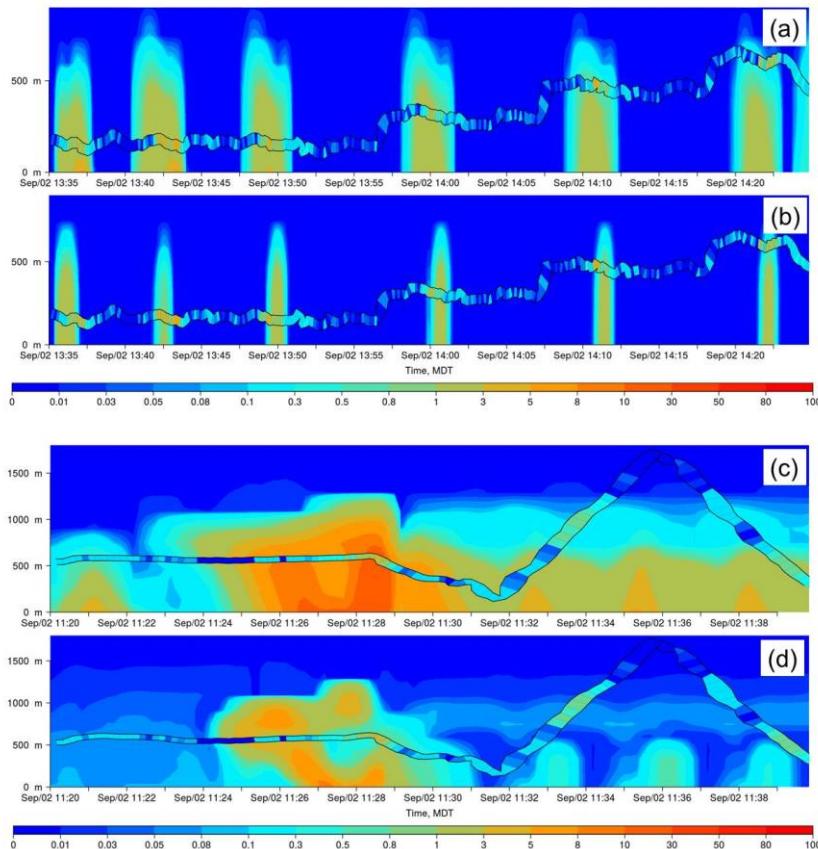
1054 These and the snap-shot comparisons described in Section 3.1 show that the higher resolution model is having a
 1055 significant impact on predictions – however, other aspects of the overall model performance are preventing the
 1056 potential benefits of higher resolution from influencing the standard performance evaluation.

1059

Table 8. Standard performance evaluation of Flight 17 for NO₂ (ppbv)

	OS2.5km	OS1km
<u>Index of Agreement</u> <u>IoA</u>	0.26	0.58
<u>Pearson Correlation Coefficient</u> <u>r</u>	0.26	0.25
<u>Normalized Mean Gross Error</u> <u>NGME</u>	2.03	1.15
<u>Mean Gross Error</u> <u>GME</u>	0.52	0.29
<u>Coefficient of Error</u> <u>CoE</u>	-0.48	0.16
<u>Root Mean Square Error</u> <u>RMSE</u>	1.37	0.70
<u>Normalized Mean Bias</u> <u>NMB</u>	0.83	-0.54
<u>Mean Bias</u> <u>MB</u>	0.21	-0.14

← Formatted Table



1060

1061
1062
1063 Figure 9. Flight 17 comparison for NO₂ (ppbv) for portions of the net flight track circling the CNRL facility for
OS2.5km (a) and OS1km (b) simulations, and for a later section of the same flight path for the OS2.5km (c) and
OS1km (d) simulations.

1064
1065

Formatted: Indent: Left: 0 cm, Hanging: 0.75 cm, Outline numbered + Level: 1 + Numbering Style: 1, 2, 3, ... + Start at: 1 + Alignment: Left + Aligned at: 0.84 cm + Indent at: 1.48 cm

1066 **4. Discussion**

1067
1068 A key result of our current work is that 1km grid cell size simulations resulted in improved prediction of plume
1069 concentration maxima relative to 2.5km grid cell size simulations, despite having no improvement using standard
1070 scoring methodologies. We also have described a scoring approach wherein these potential advantages of higher
1071 resolution may be quantified. We believe that flow field effects such as described in Figure 3 are a general result of
1072 increasing grid resolution, but note important caveats, which include:

Formatted: Justified, Indent: Left: 0 cm, Line spacing: 1.5 lines

Formatted: Justified, Line spacing: 1.5 lines

1073 (4) The availability of meteorological observation and high resolution emissions data to provide model driving
1074 information, and the resolution and proximity of this information to the simulation location. Both will
1075 influence the relative importance of grid cell size on model results. If this information is available in a
1076 higher resolution than the lower of two grid cell size simulations being compared, and/or is used via data
1077 assimilation to improve model initial meteorological conditions, our expectation is that the smaller grid cell
1078 size model may outscore the larger grid cell size model, even for more standard metrics.

1079 (5) The extent to which local, versus synoptic, weather conditions drive flow in a given region. For example, in
1080 the urban heat island meteorological simulations of Leroyer *et al.* (2014), the accuracy of local flow
1081 predictions was shown to be extremely dependent on the representation of the urban heat island, and the
1082 accuracy of the latter was critically dependent on the grid cell size (which in this example went down to 250
1083 m). In this respect, for meteorological conditions wherein local factors can dominate the flow, and where
1084 those conditions may be adequately modelled only at very high resolution, we would again expect the
1085 smaller grid cell size simulation to provide better performance, for standard metrics.

1086 (6) Conversely, model performance using standard metrics should not be expected to *increase* with
1087 successively larger and larger grid sizes; the accuracy of even the synoptic flow field will not be captured as
1088 model resolution decreases.

1089 Given these considerations, we recommend that modellers should attempt successively smaller grid cell sizes to
1090 determine the following: first, the point at which, for their particular system and simulation location, subsequent
1091 grid cell size reductions fail to improve performance; and second, to make use of still higher resolutions for studies
1092 wherein the point-to-point comparison is less important, and other factors such as accurately capturing the plume
1093 chemistry are more crucial.

1094 **45. Summary and Conclusions**

Formatted: Font: 14 pt

1095 Our work suggests the following:

1096 Decreasing ~~to~~ air-quality model horizontal grid cell size will not necessarily result in improvements to model
1097 performance in standard performance evaluations, in which the model values at the grid-cells encompassing
1098 measurement location stations are used in a pairwise comparison to observations. Other considerations, such as

1099 the accuracy of the larger scale wind direction and speed forecast, and the accuracy of the plume rise
1100 parameterization used within the model may play a greater role in the overall performance of the model, and
1101 reduce the benefits of the smaller grid cell size. In the context of a standard model performance evaluation, there
1102 may be fixed limits to the benefits of decreasing model grid cell size.

1103 Despite this difficulty, our results also show that the use of smaller grid cell sizes have some potential benefits, in
1104 that these models do a better job of resolving specific air pollution features, like high concentration maxima
1105 within plumes. Both coarse and fine grid cell size plumes may be misplaced in both time and space, with the net
1106 result that the latter model has a worse performance in a standard comparison, but is nevertheless more likely to
1107 capture the correct in-plume concentrations, and hence the chemistry, of the actual plume, in the *neighbourhood*
1108 of the observation location. When the evaluation is broadened to find the closest fit to observations in the vicinity
1109 of the observation station, with models confined to a similar representative area around the observation station,
1110 these potential benefits of the smaller grid cell size become apparent.

1111 Our results should not be taken as an indication that the standard metrics for model comparison are in some way
1112 flawed – they provide the most rigorous method for evaluating the performance of a model at specific monitoring
1113 locations and specific times. However, the ancillary performance assessment methodology presented here shows
1114 that models with very small grid sizes, which may have standard performance metric scores that have not
1115 improved or even have degraded relative to larger grid cell size models, nevertheless have scientific value, in
1116 terms of being better able to capture plume concentrations and hence plume chemistry, if not plume position.
1117 The work also suggests that the prediction accuracy of very local transport conditions may be a large factor in
1118 preventing the smaller grid cell size models from achieving improved performance in standard performance
1119 analyses.

1120 These findings suggest that at the current state of development, VHR air-quality models are of benefit for the
1121 specific purpose of chemical process studies, in which the main aim of the work is to accurately simulate plume
1122 chemistry – and in which accurate forecasting of the *position* of the plume in time and space is a secondary
1123 concern. Our work also suggests that efforts to improve other aspects of the overall modelling framework which
1124 improve the large scale flow (for example, the use of data assimilation of local meteorology to improve wind
1125 direction predictions) may result in greater benefits as smaller grid cell sizes are employed.

1126
1127 *Acknowledgements.* The authors ~~wish to~~ thank the support of Environment and Climate Change Canada
1128 (ECCC), under the CCAP program, for supporting this research. The authors also gratefully acknowledge the
1129 assistance of Michel Valin and Sylvie Gravel for advice and assistance with the installation of GEM-MACH on the
1130 Carleton University workstations during the early stages of this project.

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., Pollicastro, 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.
- Carlaw, 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/cmag/science_documentation/, last accessed September 2, 2018.
- Emery, C., Liu, Z., Russell, A.G., Talat Odman, M., Yarwood, G., and Kumar, N., Recommendations on statistics and benchmarks to assess photochemical model performance, *J. Air Waste Manage. Assoc.*, 67, 528-598, 2017.
- 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.

Garcia-Menendez, F., Hu, Y., Odman, M.T., Simulating smoke transport from wildland fires with a regional-scale air quality model : sensitivity to spatiotemporal allocation of fire emissions, *Sci. Tot. Env.*, 544-553, 2014.

Formatted: English (United States)

Formatted: Font: Italic

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

Jm, U., Bianconi, R., Solazzo, E., Kioutsioukis, I., Badia, A., Balzarini, A., Baró, R., Bellasio, R., Brunner, D., Cheme, C., Curci, G., van der Gon, H.D., Flemming, J., Forkel, R., Giordano, L., Jiménez-Guerrero, P., Hirtl, M., Hodzic, A., Honzak, L., Jorba, O., Knot, C., Makar, P.A., Manders-Groot, A., Neal, L., Perez, J.L., Pirovano, G., Pouliot, G., San Jose, R., Savage, N., Schroder, W., Sokhi, R.S., Syrakov, D., Torian, A., Tuccella, P., Wang, K., Werhahn, J., Wolke, R., Zabkar, R., Zhang, Y., Zhang, J., Hogrefe, C., Galmarini, S., Evaluation of operational on-line-coupled regional air quality models over Europe and North America in the context of AQMEII phase 2. Part II: Particulate Matter, *Atm. Environ.*, 115, 421-411, 2015.

Formatted: Indent: Left: 0 cm, Hanging: 1.27 cm, Line spacing: 1.5 lines

Formatted: Font: +Body (Calibri)

Formatted: Font: +Body (Calibri), Not Bold

Formatted: Font: +Body (Calibri)

Formatted: Font: +Body (Calibri)

Formatted: Font: +Body (Calibri)

Formatted: Font: +Body (Calibri), Not Bold

Formatted: Font: +Body (Calibri)

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

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

[Joe, D.K., Zhang, H., DeNero, S.P., Lee, H.-H., Chen, S.-H., McDonald, B.C., Harley, R.A., and Kleeman, M.J., Implementation of a high-resolution source-oriented WRF/Chem model at the Port of Oakland, *Atm. Env.*, 82, 351-363, 2014,](#)

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.

[Kain, J.S., Fritsch, J.M., A one-dimensional entraining/detraining plume model and its application in convective parameterizations. *J. Atmos. Sci.* 47, 2784-2802, 1990.](#)

[Kain, J.S., The Kain-Fritsch convective parameterization: an update. *J. Appl. Meteorol.* 43, 170-181, 2004.](#)

[Kheirbek, I., Haney, J., Douglas, S., Ito, K., Caputo, S., Jr., Matte, T., The public health benefits of reducing fine particulate matter through conversion to cleaner heating fuels in New York City, *Env. Sci. Tech.*, 48, 13573-13582, 2014.](#)

[Kheirbek, I., Haney, J., Douglas, S., Ito, K., and Matte, T., The contribution of motor vehicle emissions to ambient fine particulate matter public health impacts in New York City: a health burden assessment, *Env. Health.*, 15:89, doi: 10.1186s12940-016-0172-6, 2016.](#)

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.

[Li, J., Georgescu, M., Hyde, P., Mahalov, A., and Moutaoui, M., Achieving accurate simulations of urban impacts on ozone at high resolution, *Env. Res. Lett.*, 9, 114019 \(11pp\), 2014.](#)

[Leroyer, S., Belair, S., Husain, S.Z., and Mailhot, J., Subkilometer numerical weather predictions in an urban coastal area: a case study over the Vancouver Metropolitan Area, *J. App. Met. Clim.*, 53, 1433-1453, 2014.](#)

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 NO_x control technologies on aerosol nucleation in the source plumes, *Atm. Chem. Phys.*, 12, 11519-11531, 2012.

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,

Formatted: English (United States)

Formatted: English (United States)

Formatted: Font: Italic

Formatted: English (United States)

Formatted: Justified, Indent: Left: 0 cm, Hanging: 1.27 cm, Line spacing: 1.5 lines, No widow/orphan control

Formatted: Font color: Auto

Formatted: Font: +Body (Calibri), 11 pt, Font color: Auto

Formatted: Font: +Body (Calibri), 11 pt, Font color: Auto

Formatted: Font: +Body (Calibri), 11 pt, Font color: Auto

Formatted: Font color: Auto, English (United States)

Formatted: Font: +Body (Calibri), 11 pt, Font color: Auto

Formatted: Font color: Auto, English (United States)

Formatted: Font color: Auto, English (United States)

Formatted: Font color: Auto, English (United States)

Formatted: Font: +Body (Calibri), 11 pt, Font color: Auto

Formatted: Font: +Body (Calibri), 11 pt, Font color: Auto

Formatted: Font color: Auto

Formatted: Font: +Body (Calibri), 11 pt, Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font: Italic

- 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. Env.*, 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.

Formatted: Justified, Indent: Hanging: 1.27 cm, Line spacing: 1.5 lines, No widow/orphan control, Don't adjust space between Latin and Asian text, Don't adjust space between Asian text and numbers, Pattern: Clear

Sundqvist, Parameterization of condensation and associated clouds in models for weather prediction and general circulation simulation. In: Schlesinger M.E. (eds) *Physically-Based Modelling and Simulation of Climate and Climatic Change*. NATO ASI Series (Series C: Mathematical and Physical Sciences), vol 243. Springer, Dordrecht, 1988.

Formatted: Font: +Body (Calibri), 11 pt

Thompson, T.M., and Selin, N.E., Influence of air quality model resolution on uncertainty associated with health impacts. *Atm. Chem. Phys.*, 12, 9753-9762, 2012.

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.

5.6. Appendix A: Model Evaluation Statistics

Table A1: Model Comparison Statistics

Table A1. Model comparison statistics

Metric and Formula	Range	Ideal Score
<i>IOAI</i> ndex of Agreement (IOA) $= \begin{cases} 1 - \frac{\sum M_i - O_i }{2(\bar{O}_i - \bar{O})}, & \text{when } \sum M_i - O_i \leq 2(\bar{O}_i - \bar{O}) \\ \frac{2(\bar{O}_i - \bar{O})}{\sum M_i - O_i } - 1, & \text{when } \sum M_i - O_i > 2(\bar{O}_i - \bar{O}) \end{cases}$	[-1,1]	1
<i>CoE</i> efficient of Error (COE) $= 1 - \frac{\sum M_i - O_i }{(\bar{O}_i - \bar{O})}$	$[-\infty, 1]$	1
Mean Bias (MB) $MB = \frac{1}{N} \sum (M_i - O_i) = \bar{M} - \bar{O}$		0
Mean Gross Error (MGE) $GGE = \frac{1}{N} \sum M_i - O_i $		0
Normalized Mean Bias (NMB) $NMB = \frac{\sum (M_i - O_i)}{\sum O_i} = \left(\frac{\bar{M}}{\bar{O}} - 1 \right)$		0
Normalized Mean Gross Error (NMGE) $MGE = \frac{\sum M_i - O_i }{\sum O_i}$		0
Root Mean Square Error (RMSE) $RMSE = \sqrt{\frac{1}{N} \sum (M_i - O_i)^2}$		0
Pearson Correlation Coefficient (r) $r = \frac{\sum (M_i - \bar{M})(O_i - \bar{O})}{\sqrt{\sum (M_i - \bar{M})^2 \sum (O_i - \bar{O})^2}}$	[-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 SO₂ 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
<u>NGMEMGE</u>	1.739	1.982	1.798	0.799	0.660	0.356
<u>GMEGMGE</u>	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 SO₂ 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
<u>NGMEMGE</u>	3.143	3.281	3.215	1.896	1.300	0.964
<u>GMEGMGE</u>	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
<u>NGMEMGE</u>	0.927	1.009	0.962	0.650	0.519	0.380
<u>GMEGMGE</u>	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

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
<u>NGMEMGE</u>	1.913	1.982	1.971	1.366	0.920	0.777
<u>GMEGME</u>	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 O₃ 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
<u>NGMEMGE</u>	0.660	0.621	0.627	0.508	0.464	0.361
<u>GMEGME</u>	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 O₃ 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
<u>NGMEMGE</u>	0.706	0.775	0.773	0.600	0.361	0.352
<u>GMEGME</u>	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

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

	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
NMEMGE	0.816	0.837	0.827	0.657	0.579	0.487
GMEGME	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) ($\mu\text{g m}^{-3}$)

	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
NMEMGE	0.782	0.804	0.801	0.642	0.512	0.457
GMEGME	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