

---

# Averaging Atmospheric Gas Concentration Data using Wasserstein Barycenters

---

**Mathieu Barré**  
INRIA & CS Dept.  
Ecole Normale Supérieure,  
PSL Research University  
Paris, France.

**Clément Giron, Matthieu Mazzolini**  
Kayrros SAS  
Paris, France.

**Alexandre d'Aspremont**  
CNRS & CS Dept.  
Ecole Normale Supérieure,  
PSL Research University  
Paris, France.

## Abstract

Hyperspectral satellite images report greenhouse gas concentrations worldwide on a daily basis. While taking simple averages of these images over time produces a rough estimate of relative emission rates, atmospheric transport means that simple averages fail to pinpoint the source of these emissions. We propose using Wasserstein barycenters coupled with weather data to average gas concentration data sets and better concentrate the mass around significant sources.

## 1 Introduction

Thanks to lower launch costs and a renewed focus on earth observation, there are now several constellations of satellites monitoring greenhouse gas emissions from sun-synchronous orbits. These satellites, notably Sentinel-5P from the European Union's Copernicus program, provide daily hyperspectral images of the entire globe at a resolution of  $5.5 \times 7$  km. While anthropogenic emissions of carbon dioxide are drowned by natural sources in short time windows, this is not the case for at least two other important gases, methane ( $\text{CH}_4$ ) a very potent greenhouse gas, and nitrogen dioxide ( $\text{NO}_2$ ), a pollutant. We focus on methane here, as it has a longer lifespan than  $\text{NO}_2$  (12 years versus a few hours) hence simple averages of concentration images seriously lack contrast.

Anthropogenic methane emissions come from mostly two sources: the oil and gas industry on one hand, agriculture and waste management on the other. Methane has a lifespan of 12 years, much shorter than  $\text{CO}_2$ , but a global warming potential (GWP) that is 85 five times higher than that of  $\text{CO}_2$ , over a 20 years window. This means that mitigating methane emissions can have a very significant short term impact on warming. The latest methane budget [Saunois et al., 2020] estimates anthropogenic methane emissions at about 380 Tg  $\text{CH}_4/\text{y}$  (bottom-up), of which 108 Tg  $\text{CH}_4/\text{y}$  is coming from the oil and gas sector and 227 Tg  $\text{CH}_4/\text{y}$  are attributed to agriculture and waste management. Oil and gas emissions are concentrated in a few dense clusters around shale basins such as the Permian in the US, pipelines and major fields in e.g. Turkmenistan or Algeria. Countries and companies with high operational standards tend to have lower emission rates, and the fact that the list of key emitters is relatively short means that the oil and gas sector is both a low cost and high short term impact greenhouse gas emissions mitigation target.

Here, we solve the Wasserstein barycenter problem defined in [Rabin et al., 2011, Agueh and Carlier, 2011] on emission data sets. Recent advances in computational optimal transport using simple iterative methods [Sinkhorn, 1964, Wilson, 1969, Cuturi, 2013, Chizat et al., 2018a] mean that this task is now feasible at scale on satellite image data sets and we refer the reader to [Peyré and Cuturi, 2019] for a complete treatment. Using proper averages of concentration images has the power to remove part of the noise and better highlight key emitters while lowering the detection threshold.

Atmospheric transport is modeled by two main components, advection and diffusion. While most optimal transport problems usually come with little structural information on the transport model [Peyré and Cuturi, 2019], in the case of atmospheric transport we have historical weather records describing wind speed and direction, and some information on turbulence. We can use this information to remove biases in the optimal transportation problems underlying the barycenter computation. We thus adapt the local metric in transportation problems to account for known biases introduced by wind, and handle missing pixels using unbalanced transportation problems.

While getting accurate averaged concentrations is important, our main focus is in fact on attribution. We thus seek to solve an *inverse transportation problem*: given partial information on the transportation plan (from weather), we seek to identify emission source locations and flow rates consistent with observations, without solving a much heavier (and data intensive) full inversion problem. Early numerical experiments on Sentinel-5P images, reported below, show that Wasserstein averages are much better spatially correlated with oil & gas activity than simple averages.

**Notations** In what follows, we write  $h(\cdot)$  the discrete entropy. For  $x, y \in \mathbb{R}_+^n$ ,  $h(x) = \sum_{i=1}^n x_i \log(x_i)$  (with  $0 \log(0) = 0$ ) and we write  $\text{KL}(\cdot|\cdot)$  the KL divergence, with  $\text{KL}(x|y) = \sum_{i \in \text{Supp}(y)} x_i \log(x_i/y_i)$ .

## 2 Wasserstein Barycenter of Atmospheric Gas Concentrations

Consider we are given  $N \in \mathbb{N}$  squared images  $(g^{(k)})_{k \in [1, N]}$  with each  $g^{(k)} \in \mathbb{R}_+^{n \times n}$  representing gas concentration over time on the same region of earth. Pixels represent the mean methane concentration on small squared regions with same side length. In practice images do not need to be square, but we use this convention here for simplicity.

We use flattened version of the 2D gas images obtained by running through  $g^{(k)}$  row by row. Depending on the context we will refer to  $g^{(k)}$  as a 2D image in  $\mathbb{R}^{n \times n}$  or as a 1D vector in  $\mathbb{R}^{n^2}$ .

### 2.1 Problem formulation

We focus on the 1D discrete transport problem with entropic regularization [Cuturi, 2013]. Consider the transport problem

$$W^C(\mu, \nu) = \min_{P \mathbf{1} = \mu, P^T \mathbf{1} = \nu} \text{Tr}(C^T P) + \lambda h(P) \quad (1)$$

in the variable  $P \in \mathbb{R}_+^{n^2 \times n^2}$ , given non negative distribution variables  $\mu, \nu \in \mathbb{R}_+^{n^2}$ , a cost matrix  $C \in \mathbb{R}_+^{n^2 \times n^2}$  and a regularization parameter  $\lambda$ . The optimal value of this problem measures the effort to move the gas from a state  $\mu$  to a state  $\nu$ . With this formulation the total masses of  $\mu$  and  $\nu$  have to match. This is not compatible with our problem since methane particles are emitted over time (e.g by leaks), and some are disappearing from our measurements (e.g diffusion below the detection threshold, exiting the studied zone, etc.) and some pixels are often missing (because of water, clouds, etc.).

To compare distribution with different masses, one simple idea used in [Gramfort et al., 2015] is to add a new dimension (i.e dummy pixels) to the vector  $\mu$  and  $\nu$  that will contain the excess of mass of  $\mu$  (assuming  $\mathbf{1}^T \mu \geq \mathbf{1}^T \nu$ ), then apply the classical transport problem (1) to these augmented distributions. Optimal transport between distributions with different masses is often called unbalanced transport. Unbalanced transport has a nice fluid dynamics interpretation with the introduction of a source term that deals with the creation and deletion of matter during the transport [Chizat et al., 2018b, Liero et al., 2018]. In practice here, we use a formulation of unbalanced transportation that relaxes marginal constraints in (1). Given another regularization parameter  $\lambda_u$ , the problem thus

becomes

$$W_u^C(\mu, \nu) = \min_{P \in \mathbb{R}_+^{n^2 \times n^2}} \text{Tr}(C^T P) + \lambda h(P) + \lambda_u \text{KL}(P1|\mu) + \lambda_u \text{KL}(P^T 1|\nu) \quad (2)$$

We are interested in obtaining an aggregation of a sequence  $(g^{(k)})$  of methane distributions recorded daily over some period of time. The transportation cost  $W_u^C(\cdot, \cdot)$  can be used as a metric between the gas distributions  $(g^{(k)})_{k \in [1, N]}$  and, given a sequence of cost matrix  $(C^{(k)})_{k \in [1, N]}$ , the Wasserstein barycenter (WB) of the  $(g^{(k)})$  is computed as

$$\bar{g} = \operatorname{argmin} \sum_{i=k}^N W_u^{C^{(k)}}(g^{(k)}, g) \quad (\text{WB})$$

in the variable  $g \in \mathbb{R}_+^{n^2}$ . Algorithms to solve this convex optimization problem are described in e.g. [Benamou et al., 2015, Chizat et al., 2018a].

## 2.2 Choice of local cost function

The cost matrix  $C$  plays an important role in the transport problem. In our setting  $C_{ij}$  represents the effort necessary to move a particle from position  $X_j = (\lceil \frac{j}{n} \rceil, \lfloor j \bmod n \rfloor + 1)$  to position  $X_i = (\lceil \frac{i}{n} \rceil, \lfloor i \bmod n \rfloor + 1)$ .

In our experiments we use several types of cost matrix, with  $C_{ij}^{(k)} = c(X_i, X_j)$  and  $c(\cdot, \cdot)$  defined as

$$\text{Euclidean metric : } c(x, y) = \|x - y\|^2 \quad (L_2)$$

$$\text{Wasserstein Fisher Rao metric : } \delta > 0, c(x, y) = -\log \left( \cos^2 \left( \frac{\|x-y\|}{2\delta} \wedge \frac{\pi}{2} \right) \right) \quad (\text{WFR})$$

$$\text{Euclidean + Wind : } t > 0, c(x, y) = (\|x - y\|^2 - t \langle w_k, x - y \rangle)_+ \quad (L_2 + W)$$

The choice of (WFR) is motivated by [Chizat et al., 2018b, Corollary 5.9] which links the formulation of (2) with a fluid mechanics interpretation of the transport, with creation and deletion of mass. When this cost is used,  $\lambda_u$  is set to 1 and the parameter  $\delta$  is used to control mass creation.

In the case of the  $(L_2 + W)$  metric,  $w_k \in \mathbb{R}^2$  is a mean wind vector associated with image  $g^{(k)}$ . The simple idea behind this cost matrix is that, in the presence of wind, it is easier for a particle to move with the wind than against it. Indeed if the shifting vector  $x - y$  from  $y$  to  $x$  follows the direction of the wind  $w_k$  the cost is smaller, and  $t$  is a scaling parameter to balance the two terms.

## 3 Experiments

We used the optimal transport package POT [Flamary and Courty, 2017] to compute the Wasserstein barycenters in our experiments, slightly modified to support the use of multiple cost matrices.

### 3.1 Synthetic Experiments

In this section we present synthetic experiments to illustrate what can be obtained using Wasserstein barycenters and to emphasize the importance of cost matrices choices. Further intuitive examples of Wasserstein barycenters can be found in Cuturi and Doucet [2014].

First, we look at a simple setting where arithmetic mean and Wasserstein barycenter give very different results presented in Figure 1. The  $g^{(k)}$  are taken to be Gaussian clouds with unit variance and a mean  $\mu_k \in \mathbb{R}^2$  that rotates around the center of the image with  $k$  (to mimic wind changing). We observe that the arithmetic mean described a circle around the center, whereas the Wasserstein barycenter is concentrated on the center. See Appendix B for a synthetic example with  $(L_2 + W)$  cost.

### 3.2 Experiments on Real Data

The data represents  $\text{CH}_4$  concentration over a period of time from January 1 2019 until June 1 2020 (see § A for details). Experiments on the Permian basin are also presented in Appendix D.

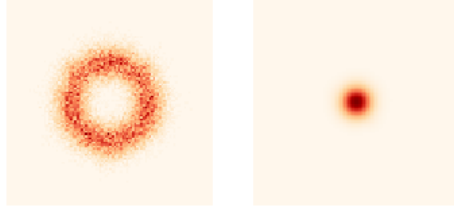


Figure 1: Left: arithmetic mean of  $g^{(k)}$ . Right: Wasserstein barycenter of  $g^{(k)}$ .

### 3.2.1 Ohio - West Virginia - Pennsylvania Mines

The first zone that we look at is a mining area located at the border between Ohio, West Virginia and Pennsylvania. In this region, wind introduces an important bias in the transport of the gas (see Appendix C for more details).

Figure 2 shows the result of barycenter computations with different cost matrices. Due to the bias West-East in the wind distribution, we see that the barycenter with costs  $(L_2 + W)$  is shifted to the right compared to barycenter with costs  $(L_2)$ . The blue dots correspond to every coal seams that have been exploited at some point in time in the region (see Appendix A.3 for sources). While Figure 2 only displays coal seams, unconventional gas wells in Southern Pennsylvania are located in the same area [Barkley et al., 2019]. Our Wasserstein barycenter with cost  $(L_2 + W)$  highlights potential methane emitters in Southwestern Pennsylvania. This particular region is well-known for being a significant source of anthropogenic methane [Barkley et al., 2019], due to underground coal and unconventional gas extraction activities. It appears that Wasserstein barycenters are much more correlated with fossil fuels production than classical means.

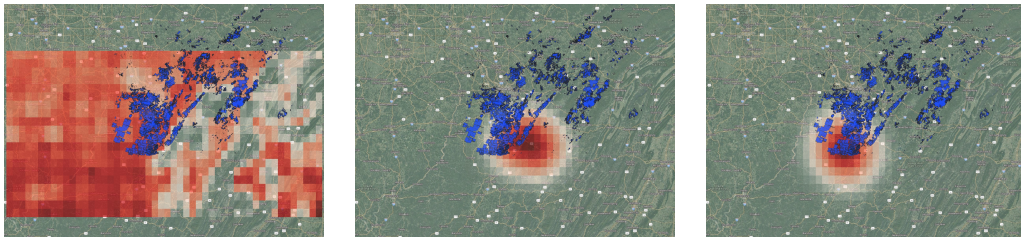


Figure 2: Left: Arithmetic mean. Middle:  $(WB) + (L_2)$ . Right:  $(WB) + (L_2 + W)$

### 3.2.2 Irak-Koweit region

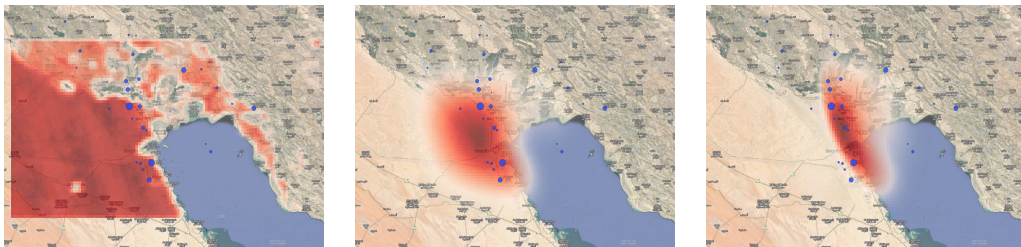


Figure 3: Left: Arithmetic mean. Middle:  $(WB) + (L_2)$ . Right:  $(WB) + (WFR)$

We also studied the region around the borders between Iraq, Iran and Kuwait, an area known for its high Oil production activity. Figure 3 displays the results of our computations. As in the previous example, blue dots correspond the oil fields in the region and the size of the dots represents production level (source: Kayrros analysis). One observes on the averaged image that the satellites measured a high concentration on all the desert areas. This might be due to the well-known albedo-induced bias in the Sentinel-5 Precursor Level 2 Methane data ESA [2020]. Using the Wasserstein-Fisher-Rao

cost function leads to a clear accumulation of the mass of emissions around the south eastern part of the country, home to the biggest oil field of the country (Rumaila, approximately 1.4 million barrels per day). Due to the high resolution of the image, we couldn't use the  $(L_2 + W)$  metric as it requires too much computation at this point.

## Acknowledgements

The authors would like to thank Thomas Lauvaux for sharing his expertise on methane production in the Ohio - West Virginia - Pennsylvania region. AA is at CNRS, and CS Department, Ecole Normale Supérieure, PSL Research University, 45 rue d'Ulm, 75005, Paris and Kayros SAS. AA would like to acknowledge support from the *ML and Optimisation* joint research initiative with the *fonds AXA pour la recherche* and Kamet Ventures, a Google focused award, as well as funding by the French government under management of Agence Nationale de la Recherche as part of the "Investissements d'avenir" program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute). Contact: [aspregon@ens.fr](mailto:aspregon@ens.fr). MB acknowledges support from an AMX fellowship.

## References

- Martial Agueh and Guillaume Carlier. Barycenters in the wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924, 2011.
- ZR Barkley, T Lauvaux, KJ Davis, A Deng, A Fried, P Weibring, D Richter, JG Walega, J DiGangi, SH Ehrman, et al. Estimating methane emissions from underground coal and natural gas production in southwestern pennsylvania. *Geophysical Research Letters*, 46(8):4531–4540, 2019.
- Jean-David Benamou, Guillaume Carlier, Marco Cuturi, Luca Nenna, and Gabriel Peyré. Iterative bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015.
- Lénaïc Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard. Scaling algorithms for unbalanced optimal transport problems. *Mathematics of Computation*, 87(314):2563–2609, 2018a.
- Lénaïc Chizat, Gabriel Peyré, Bernhard Schmitzer, and François-Xavier Vialard. Unbalanced optimal transport: Dynamic and kantorovich formulations. *Journal of Functional Analysis*, 274(11):3090–3123, 2018b.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pages 2292–2300, 2013.
- Marco Cuturi and Arnaud Doucet. Fast computation of wasserstein barycenters. *Journal of Machine Learning Research*, 2014.
- ESA. S5p mission performance center methane product (readme), 2020. URL <https://sentinel.esa.int/documents/247904/3541451/Sentinel-5P-Methane-Product-Readme-File>.
- Rémi Flamary and Nicolas Courty. Pot python optimal transport library, 2017. URL <https://pythonot.github.io/>.
- Alexandre Gramfort, Gabriel Peyré, and Marco Cuturi. Fast optimal transport averaging of neuroimaging data. In *International Conference on Information Processing in Medical Imaging*, pages 261–272. Springer, 2015.
- Matthias Liero, Alexander Mielke, and Giuseppe Savaré. Optimal entropy-transport problems and a new hellinger–kantorovich distance between positive measures. *Inventiones mathematicae*, 211(3):969–1117, 2018.
- Gabriel Peyré and Marco Cuturi. Computational optimal transport: With applications to data science. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019.
- Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Bernot. Wasserstein barycenter and its application to texture mixing. In *International Conference on Scale Space and Variational Methods in Computer Vision*, pages 435–446. Springer, 2011.

Marielle Saunois, Ann R Stavert, Ben Poulter, Philippe Bousquet, Josep G Canadell, Robert B Jackson, Peter A Raymond, Edward J Dlugokencky, Sander Houweling, Prabir K Patra, et al. The global methane budget 2000–2017. *Earth System Science Data*, 12(3):1561–1623, 2020.

Richard Sinkhorn. A relationship between arbitrary positive matrices and doubly stochastic matrices. *The annals of mathematical statistics*, 35(2):876–879, 1964.

Alan Geoffrey Wilson. The use of entropy maximising models, in the theory of trip distribution, mode split and route split. *Journal of transport economics and policy*, pages 108–126, 1969.

## A Data Sources

### A.1 Sentinel 5P

We use total column CH<sub>4</sub> (XCH<sub>4</sub>) measurements from the spaceborne Tropospheric Monitoring Instrument (TROPOMI). TROPOMI is in polar sun-synchronous orbit and provides global mapping of atmospheric methane columns on daily overpasses at about 13:30 local solar time with  $7 \times 7$  km nadir pixel resolution ( $7 \times 5.5$  km since June 2019). The mission performance report for Sentinel-5 Precursor Level 2 Methane product [ESA \[2020\]](#) states that "the averaged bias for the comparison against 22 TCCON sites is  $-0.8\%$  and  $-0.31\%$  for the standard and bias corrected XCH<sub>4</sub> product". Note that for various reasons (body of water, cloud cover, etc) a significant fraction of the pixels are missing, hence any averaging method used on these images needs to properly account for missing values. Figure 4 illustrates the variability of the data coverage worldwide.

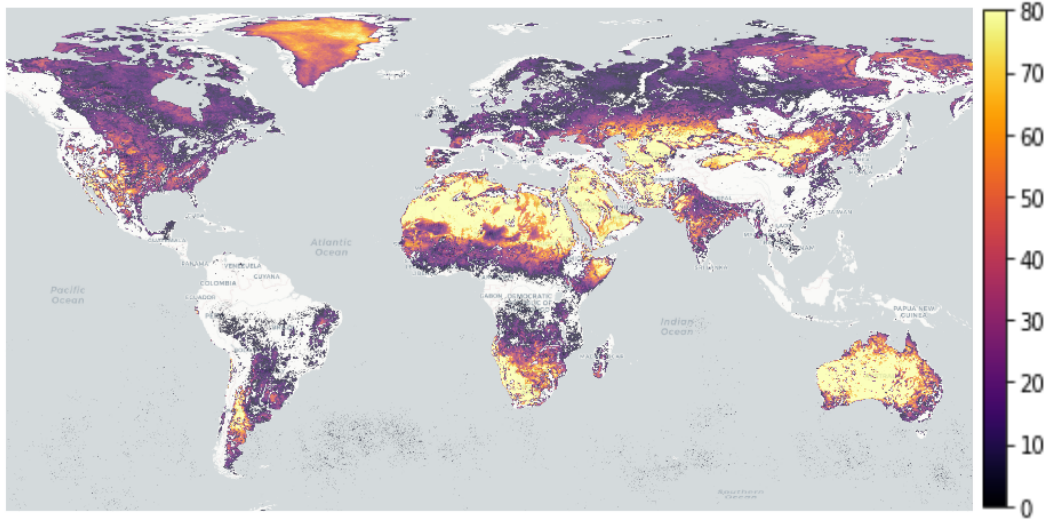


Figure 4: Sentinel-5P coverage for Level 2 XCH<sub>4</sub> data product in 2019. The value of each pixel correspond to the number of days for which Sentinel-5P provided a valid (after quality filtering, [\[ESA, 2020\]](#)) measurement for the corresponding area during year 2019. The data was clipped at 80 for clarity; some pixels exceed this value.

### A.2 Weather

Weather data are provided by the ECMWF-ERA5 product at a resolution of 0.25 degree around the equator and sample every hour. We use the northern and eastern direction of the wind measured at 100m above ground to get the direction and the intensity.

### A.3 Production Sites

- **Ohio - West Virginia - Pennsylvania:** mines locations can be found at <https://www.pasda.psu.edu/uci/DataSummary.aspx?dataset=257>.

- **Iraq-Iran-Kuweit:** Oil and Gas basins data from Kayrros
- **Permian basin:** completion data from Kayrros

#### A.4 Maps

We used google earth for the maps underlying our plots.

### B Toy Example

This toy example consists in choosing the  $g^{(k)}$  as Gaussian clouds of unit variance and mean  $\mu_k$  that drift away from the center to the right of the image. In Figure 5 we compare the results of Wasserstein barycenters with different cost matrix. For the right plot we consider the  $(L_2 + W)$  cost with a constant wind from the left to the right. Thus the barycenter has been translated in the opposite direction of the wind.



Figure 5: From left to right: arithmetic mean of the  $g^{(k)}$ , Wasserstein barycenter of the  $g^{(k)}$  with  $(L_2)$  cost, Wasserstein barycenter of the  $g^{(k)}$  with  $(L_2 + W)$  cost using constant East wind.

### C Wind in Pennsylvania

Figure 6 represents the distribution of the wind mean directions and speeds during the studied period. Bars indicates the direction from which the wind comes from (meteorological convention), the size of the bars corresponds to the directions frequencies and the color to the speed (blue is slow and red is fast). We notice that in this region the wind is mainly blowing on some West-East direction.

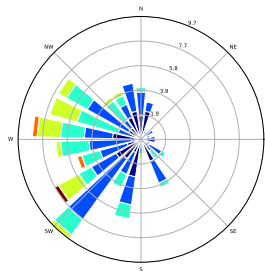


Figure 6: Distribution of the mean wind speed and direction on the Pennsylvania mining region.

### D Permian

In this section we focus on the Permian basin. This area is known for its high production of shale Oil and Gas. The result of our experiments on this region is displayed on Figure 7. The blue dots represent the locations of the recently completed oil wells (source: Kayrros analysis), over a period where production reached its maximum capacity. The dot size is proportional with the quantity of Oil and Gas produced by the well.

One can observe on Figure 7 Bottom Right that the quantity of missing observations in the satellite data is quite high. In particular, we get very few observations in a region with high source density on the East side. However Wasserstein barycenter gives some mass to this region, and using Wasserstein-Fisher-Rao metric even allows to separate the two production sub-basins. These two areas correspond to the two major oil and gas sub-basins within the Permian, namely the Delaware basin and the Midland basin.

Due to the higher resolution of these observations, we couldn't use the metric involving the wind data as it became too computationally intensive.

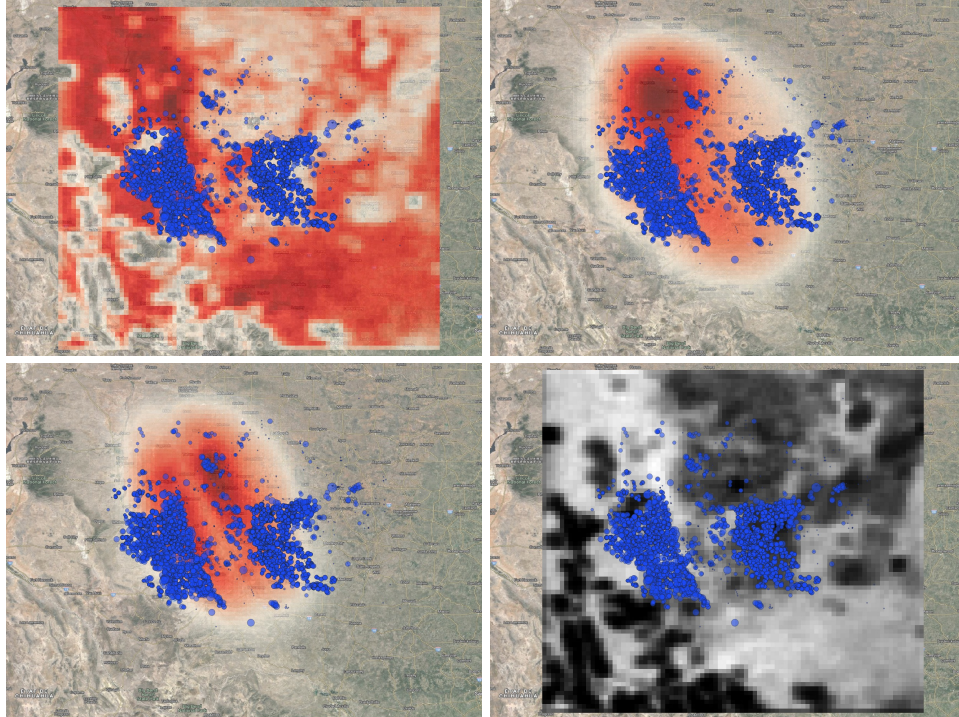


Figure 7: Top Left: Arithmetic mean. Top Right:  $(WB) + (L_2)$ . Bottom Left:  $(WB) + (WFR)$ . Bottom Right: Proportion of observed pixels (black is 0%, white is 30%).