



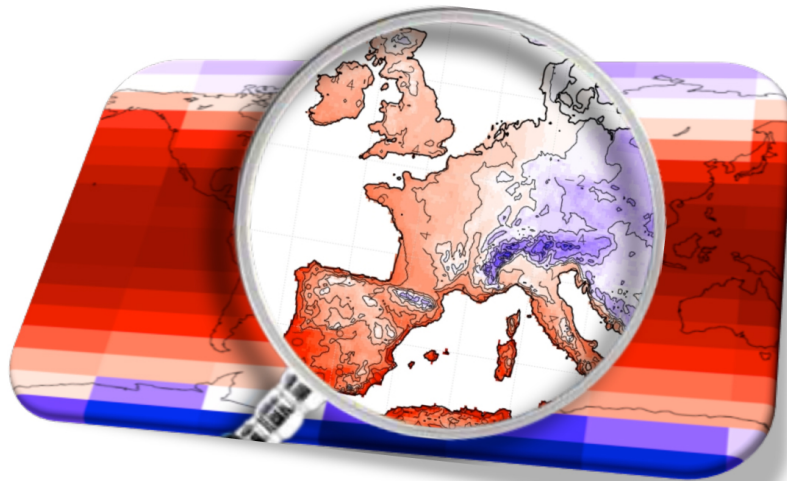
LSCE

LABORATOIRE DES SCIENCES DU CLIMAT
& DE L'ENVIRONNEMENT



Institut
Pierre
Simon
Laplace

Modélisation statistique de données climatiques : Descente d'échelle (aka downscaling) & Correction de biais - extrêmes inclus -



Mathieu Vrac

Journée SAMA (IPSL), 13 mars 2018

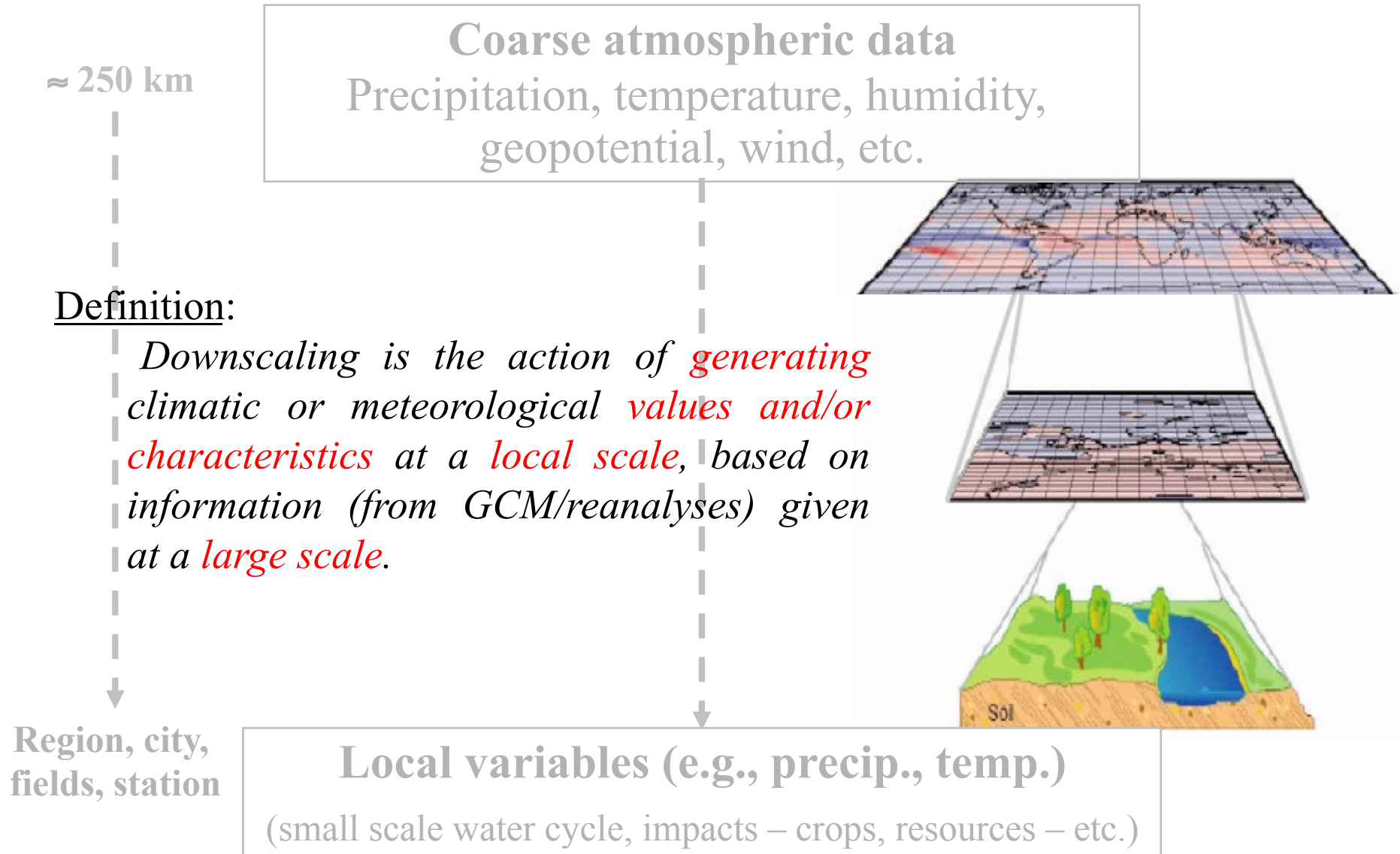
Motivations

- Many environmental/hydro/agro/human/economic activities and studies are directly affected by meteorological conditions
- **IPCC** scenarios of climate change have a **coarse spatial resolution !!**
Not adapted to the spatial scales of impact studies
 - Environmental, human, social and economic impacts
 - How will climate change interact with environmental features existing at a regional/local scale ?
 - **Downscaling**: To derive sub-grid scale (regional or local) weather or climate using General Circulation Models (GCMs) outputs or reanalysis data (e.g. NCEP)
 - **Statistical Bias Correction** also often needed !!
 - Need to improve also the **modelling of extreme events**

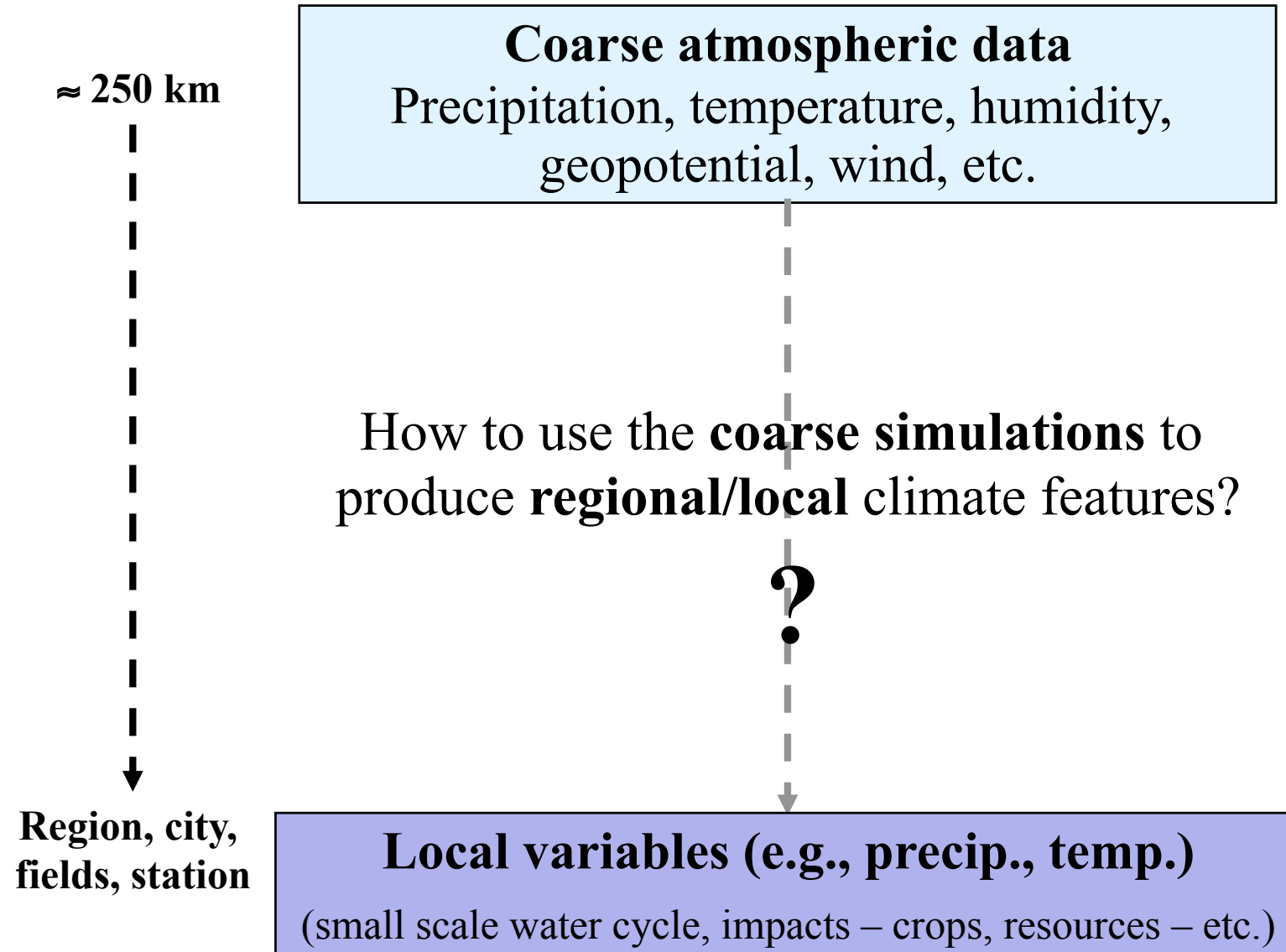
Outline

- Downscaling: main statistical approaches
(including bias correction)
- Illustration of a SWG approach
 - One extension to extremes
- Illustration of a BC approach
 - One extension to extremes
- Conclusions & perspectives

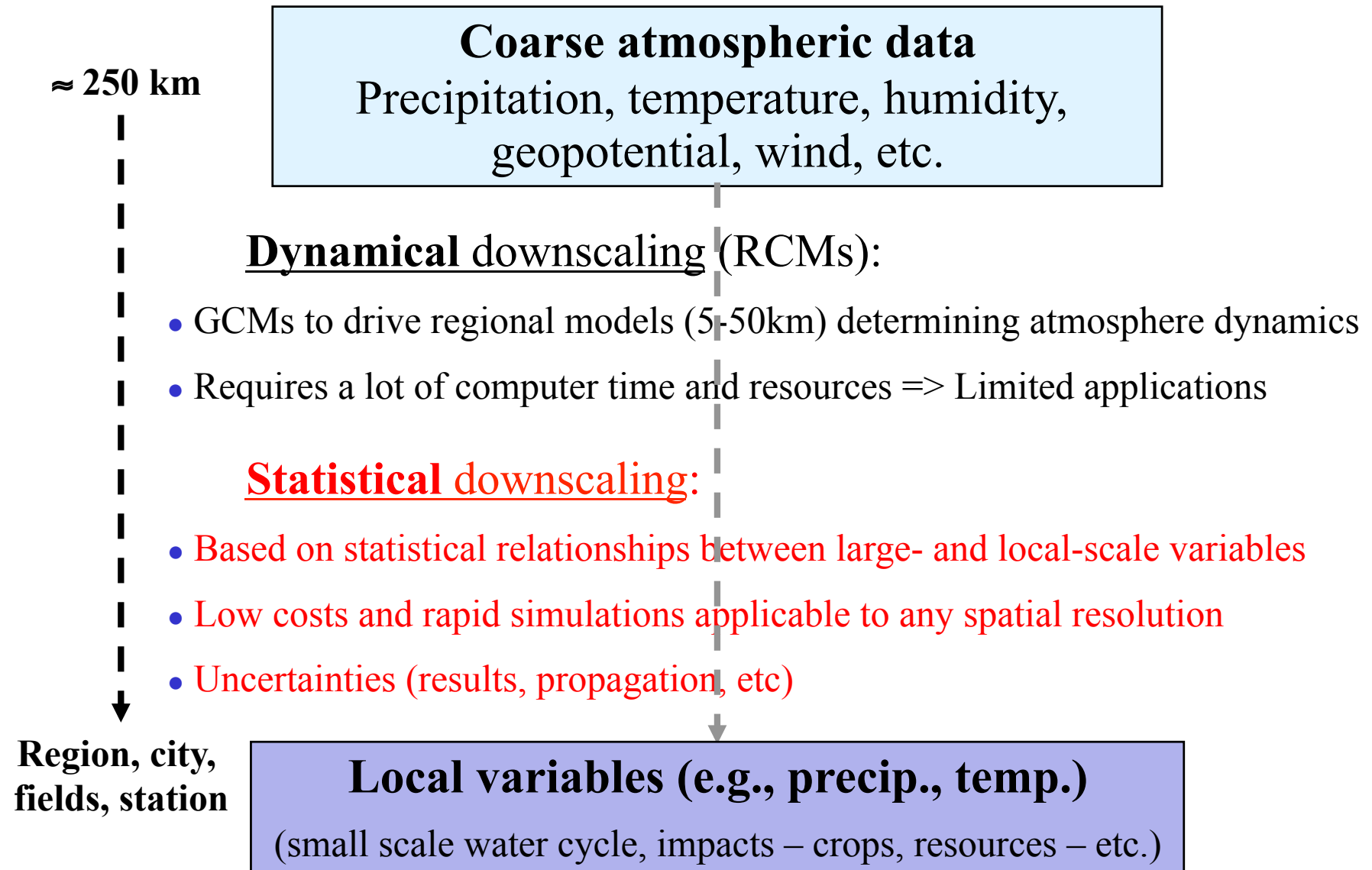
What is downscaling ???



How to downscale?: The basics



How to downscale?: The basics



Main statistical **downscaling** approaches

Predictors

$(X_i)_{i=1, \dots, p}$

Coarse atmospheric data
Precip., temp., humidity, geopot., wind, etc.

Downscaling
function

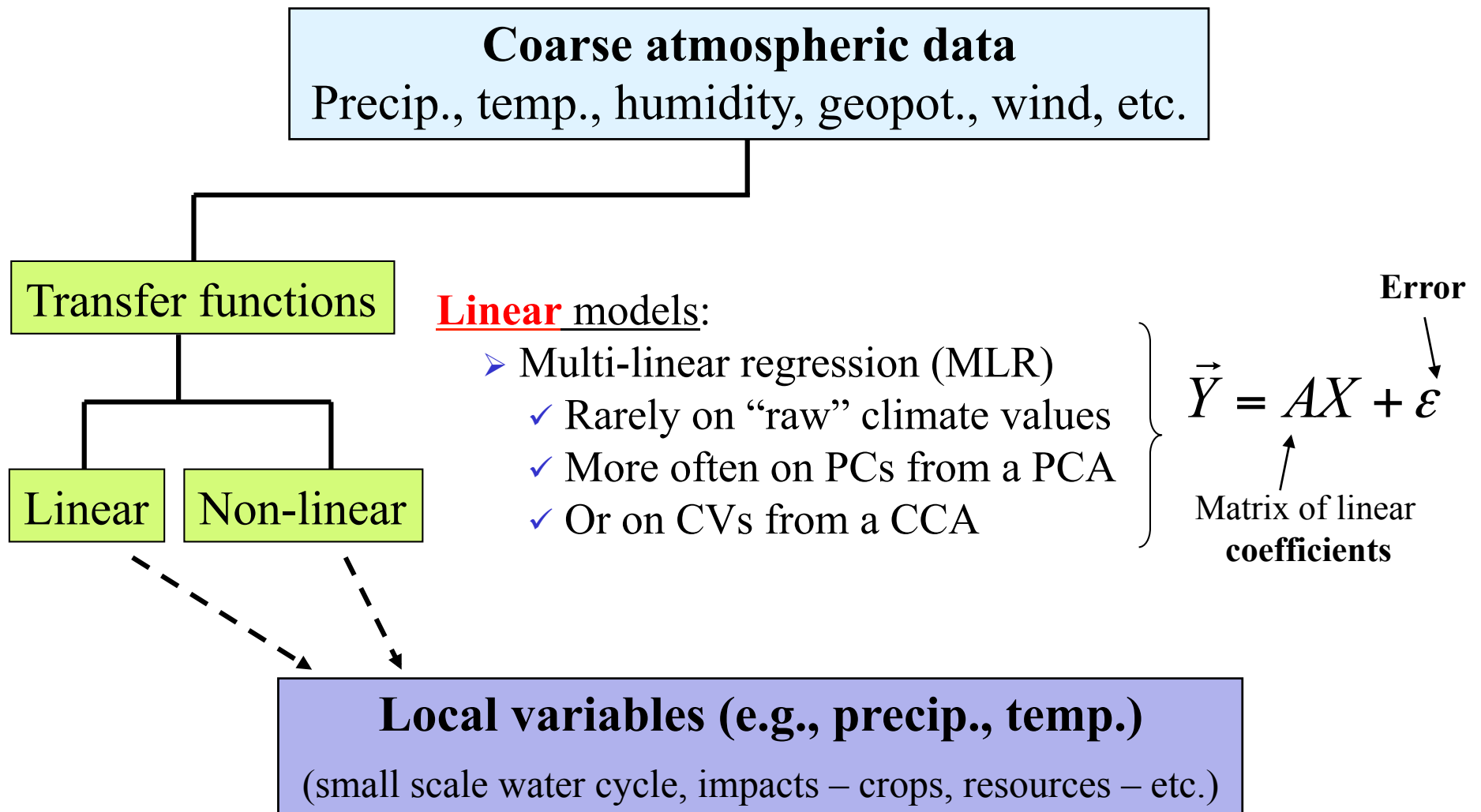
$f(X_1, \dots, X_p)$

Predictands

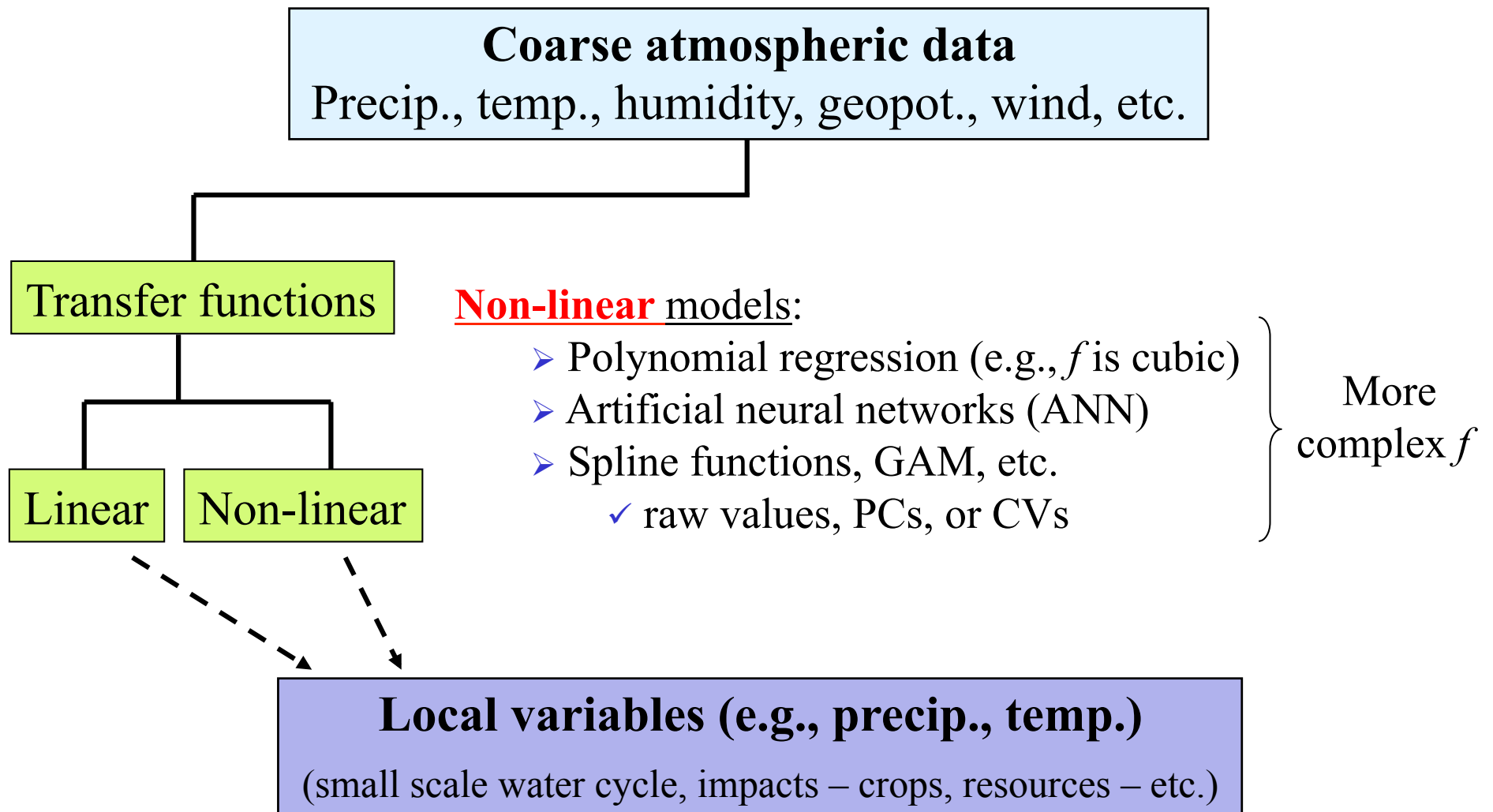
Y

Local variables (e.g., precip., temp.)
(small scale water cycle, impacts – crops, resources – etc.)

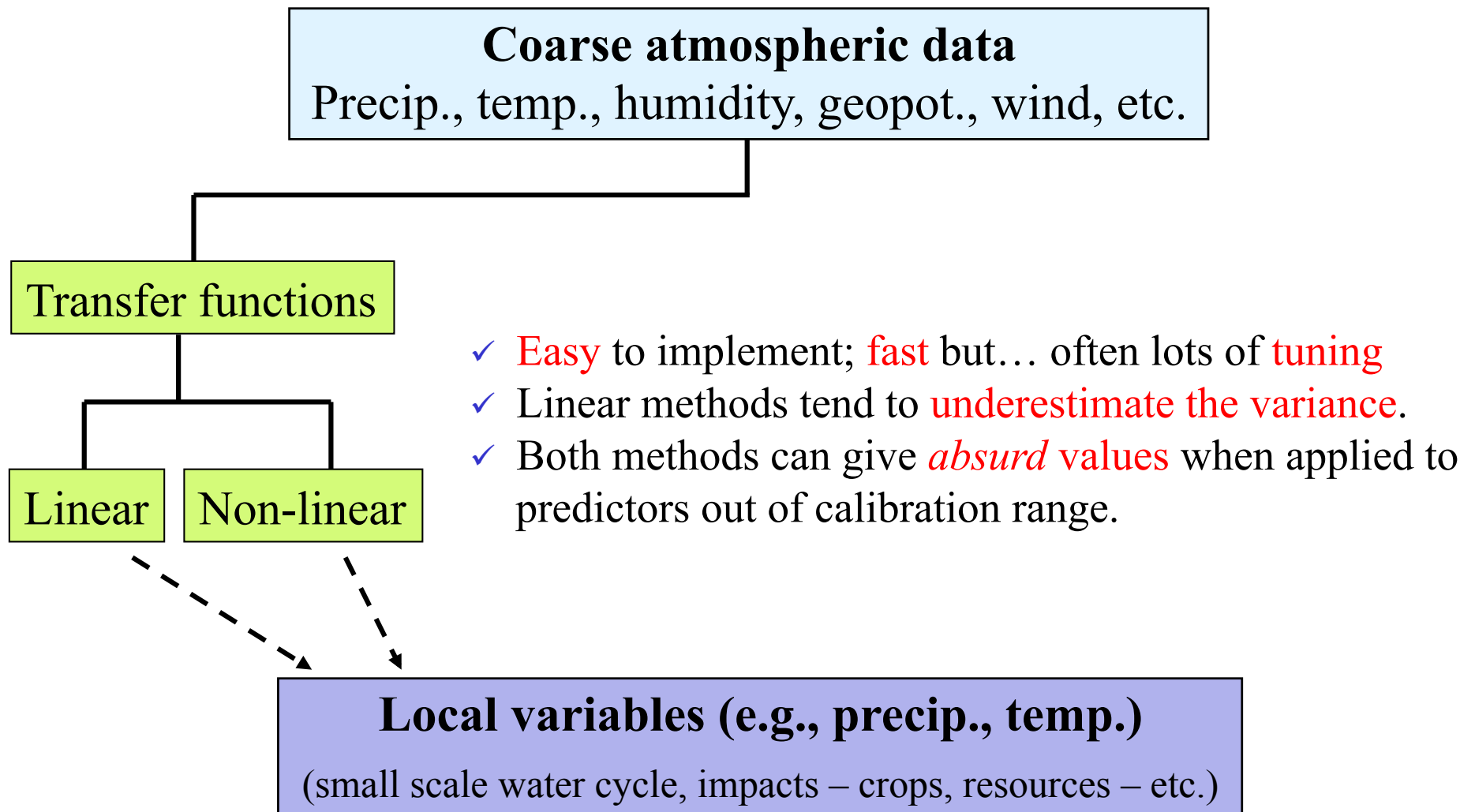
Main statistical **downscaling** approaches



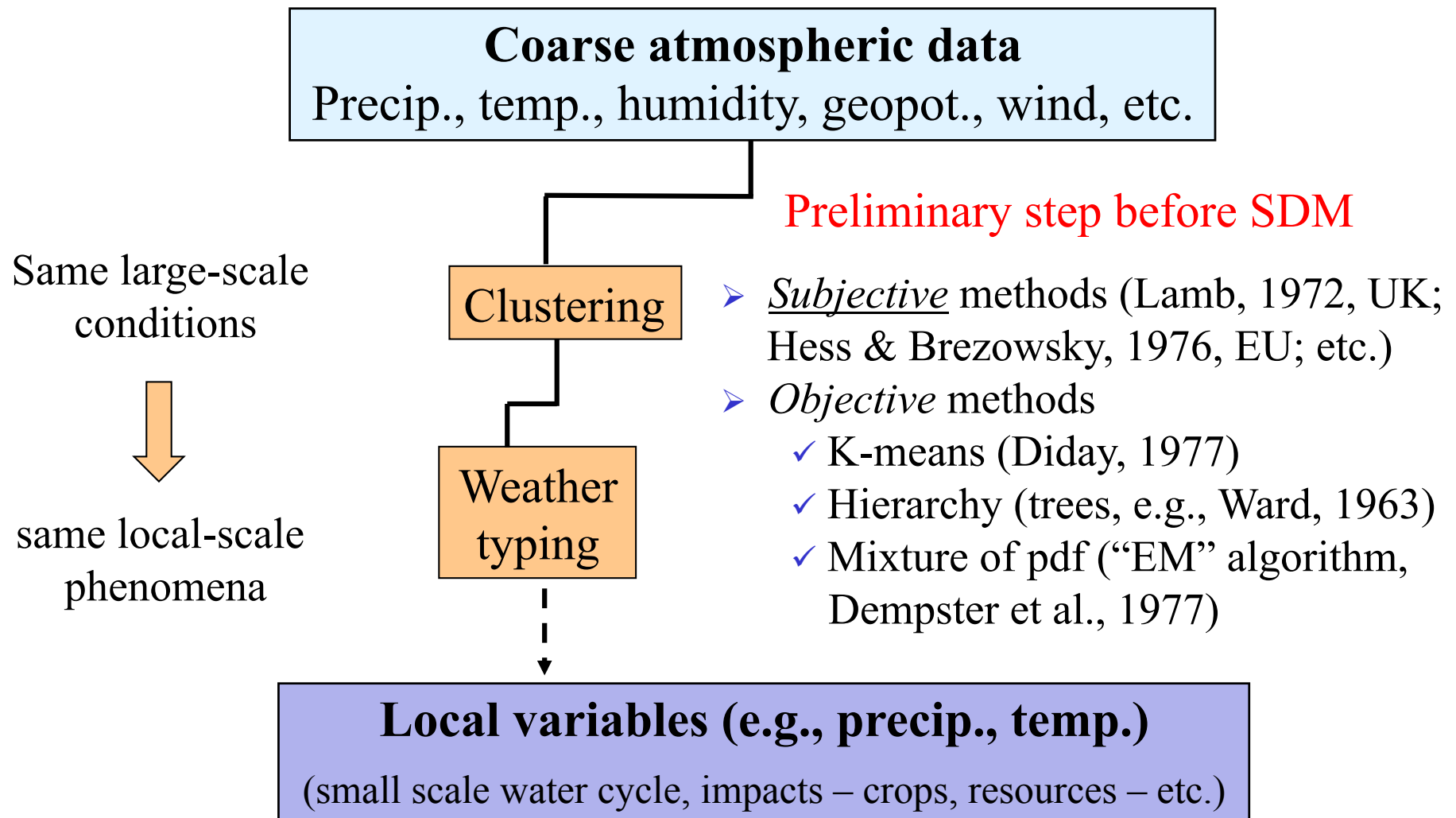
Main statistical **downscaling** approaches



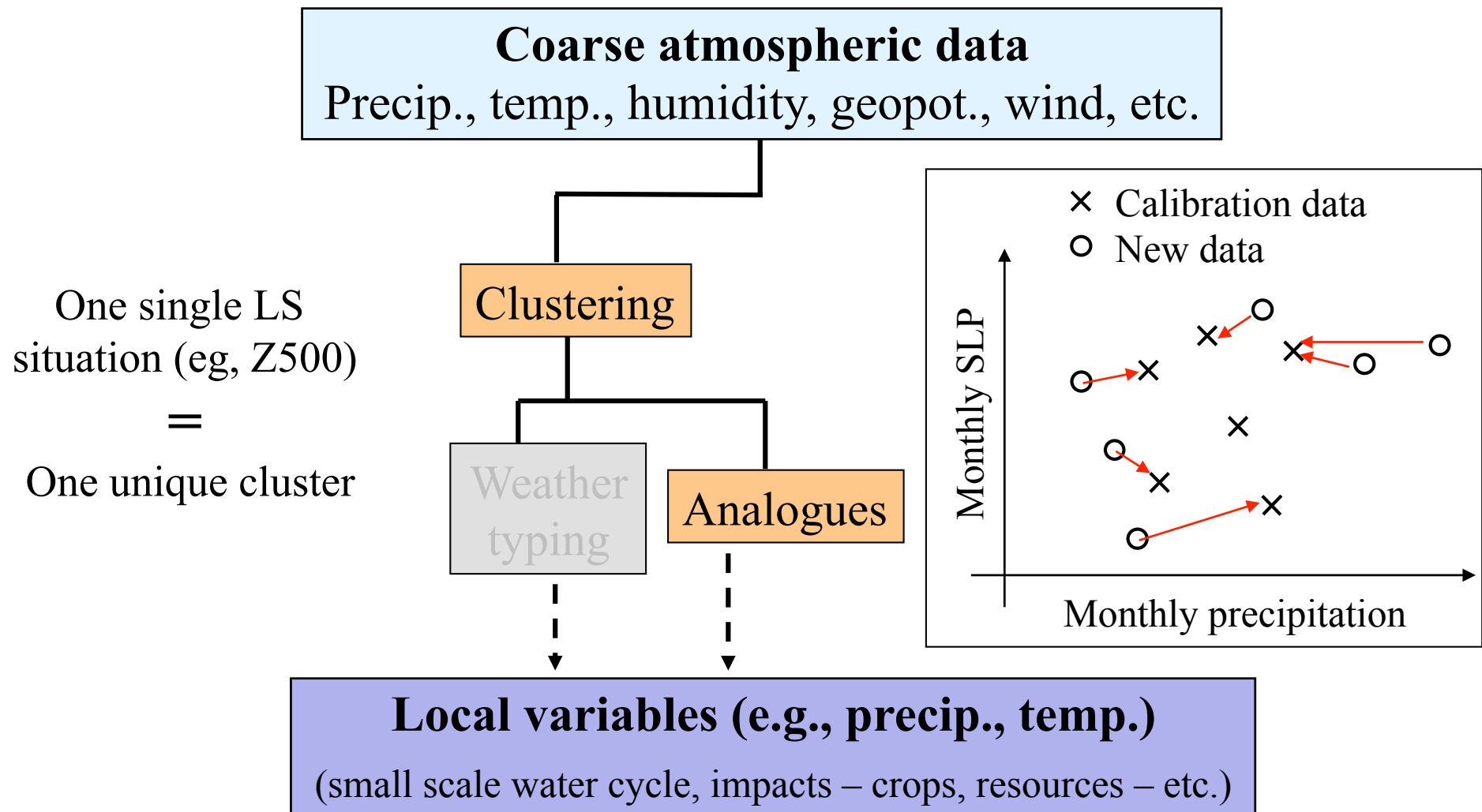
Main statistical **downscaling** approaches



Main statistical **downscaling** approaches

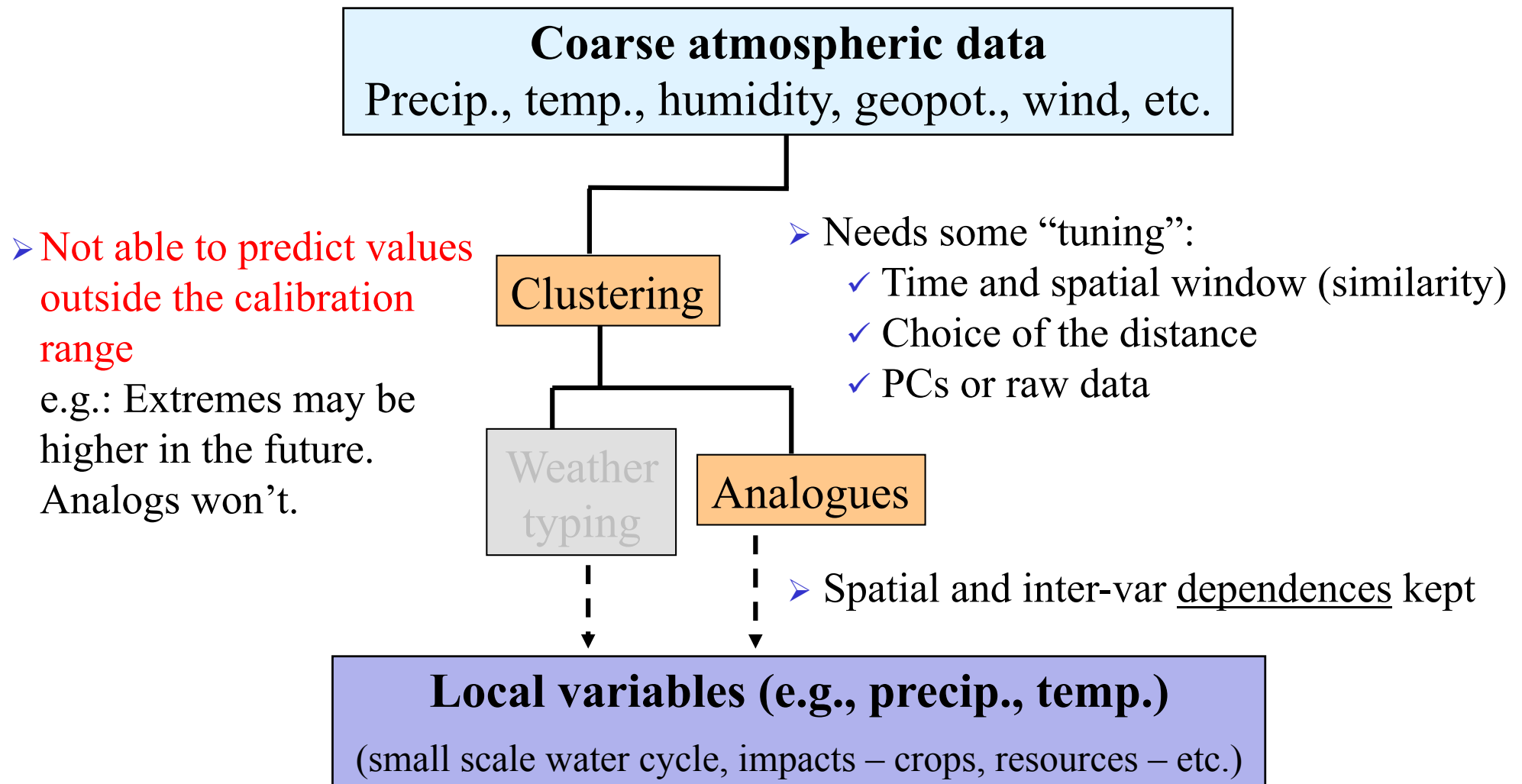


Main statistical **downscaling** approaches

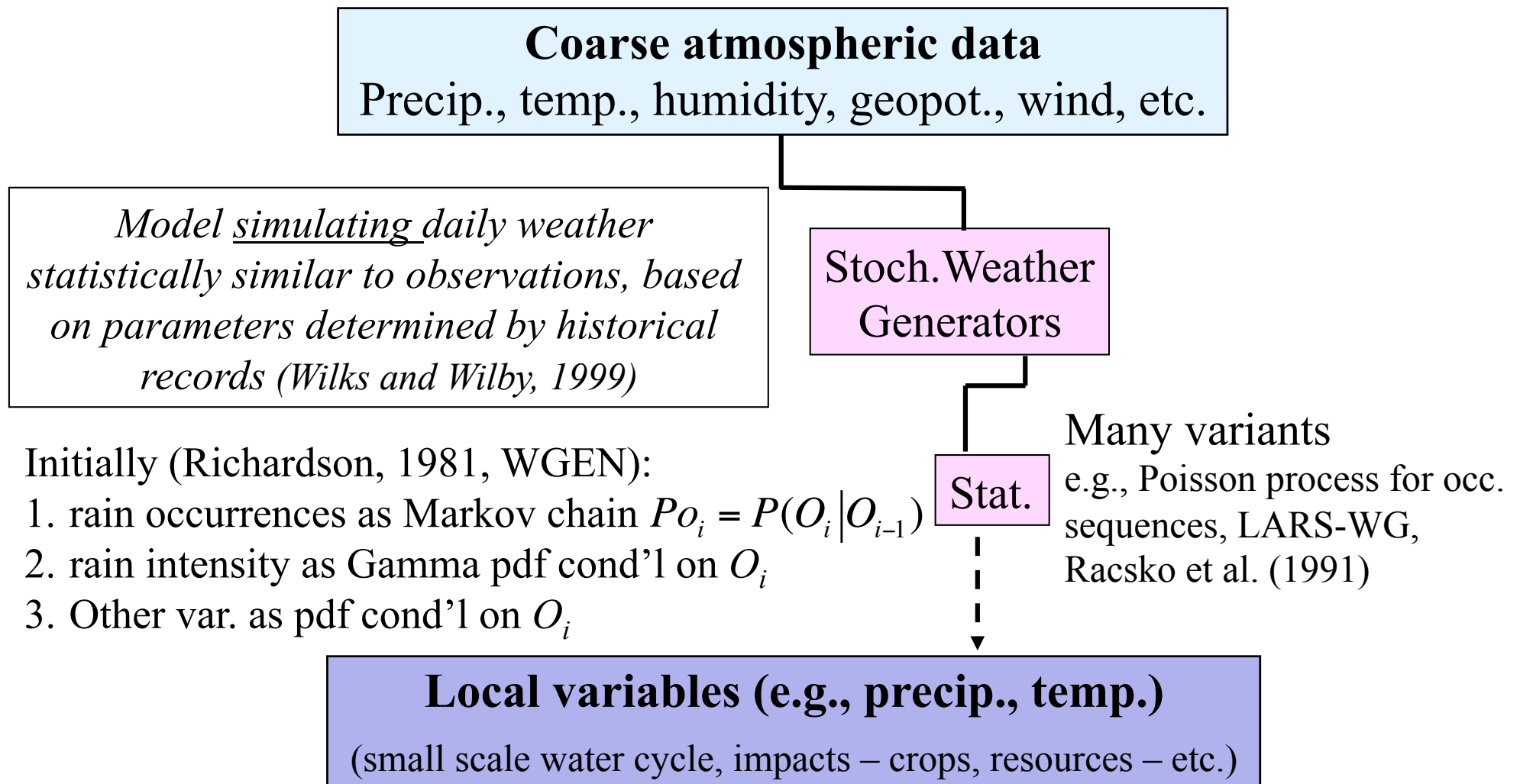


Barnett & Preisendorfer (1978); Zorita & von Storch (1998); Yiou et al. (2007, 2013,...), etc.

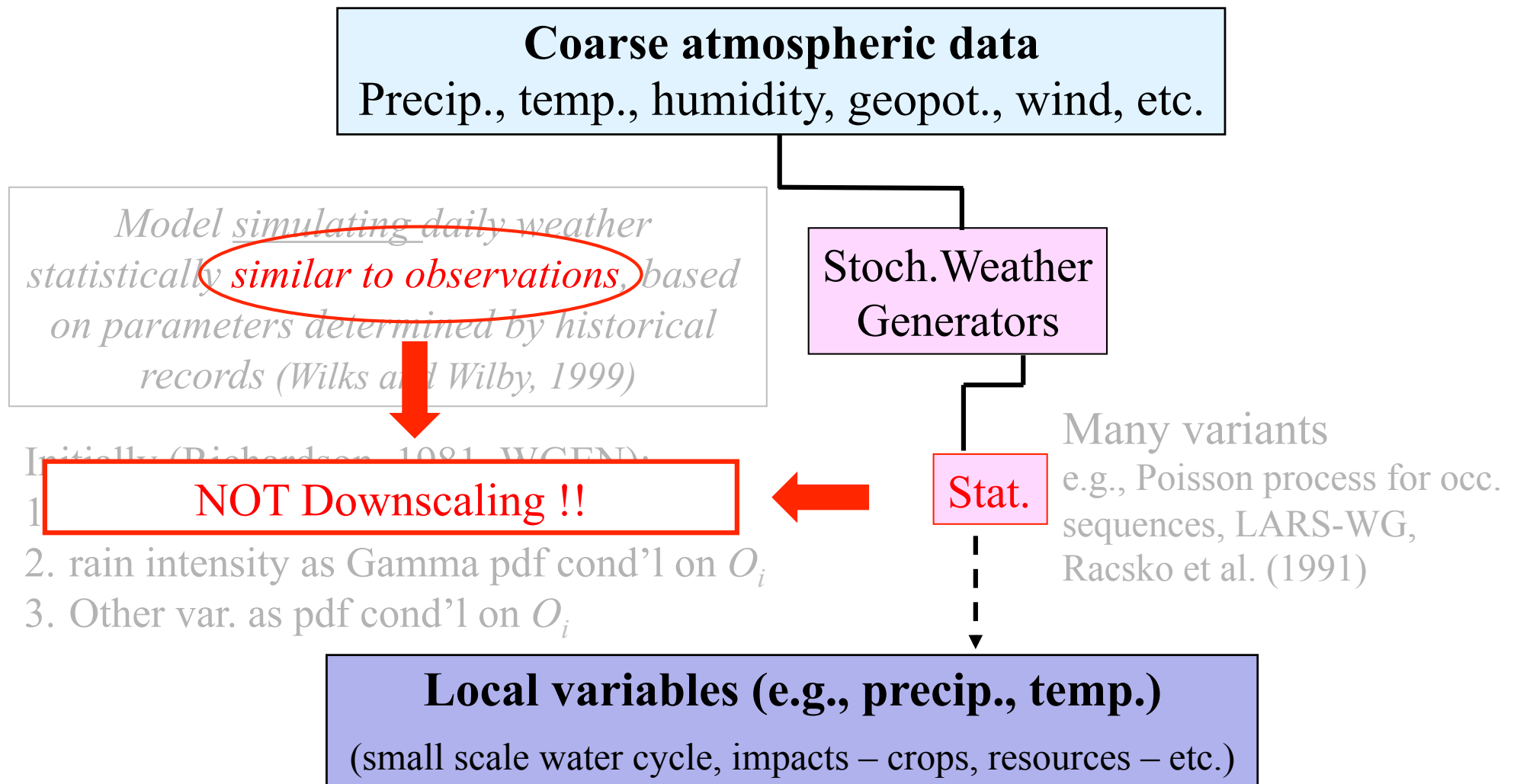
Main statistical **downscaling** approaches



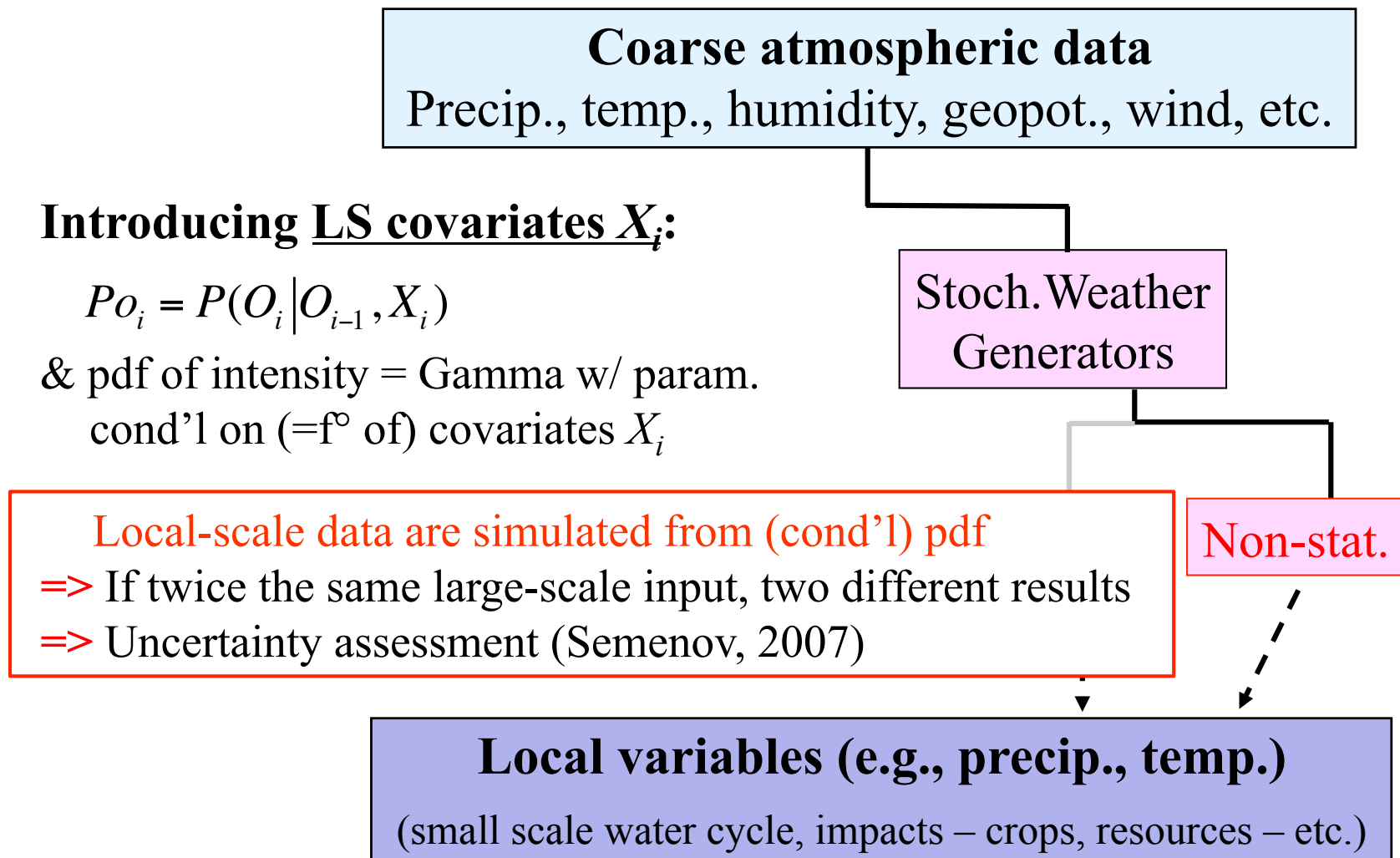
Main statistical **downscaling** approaches



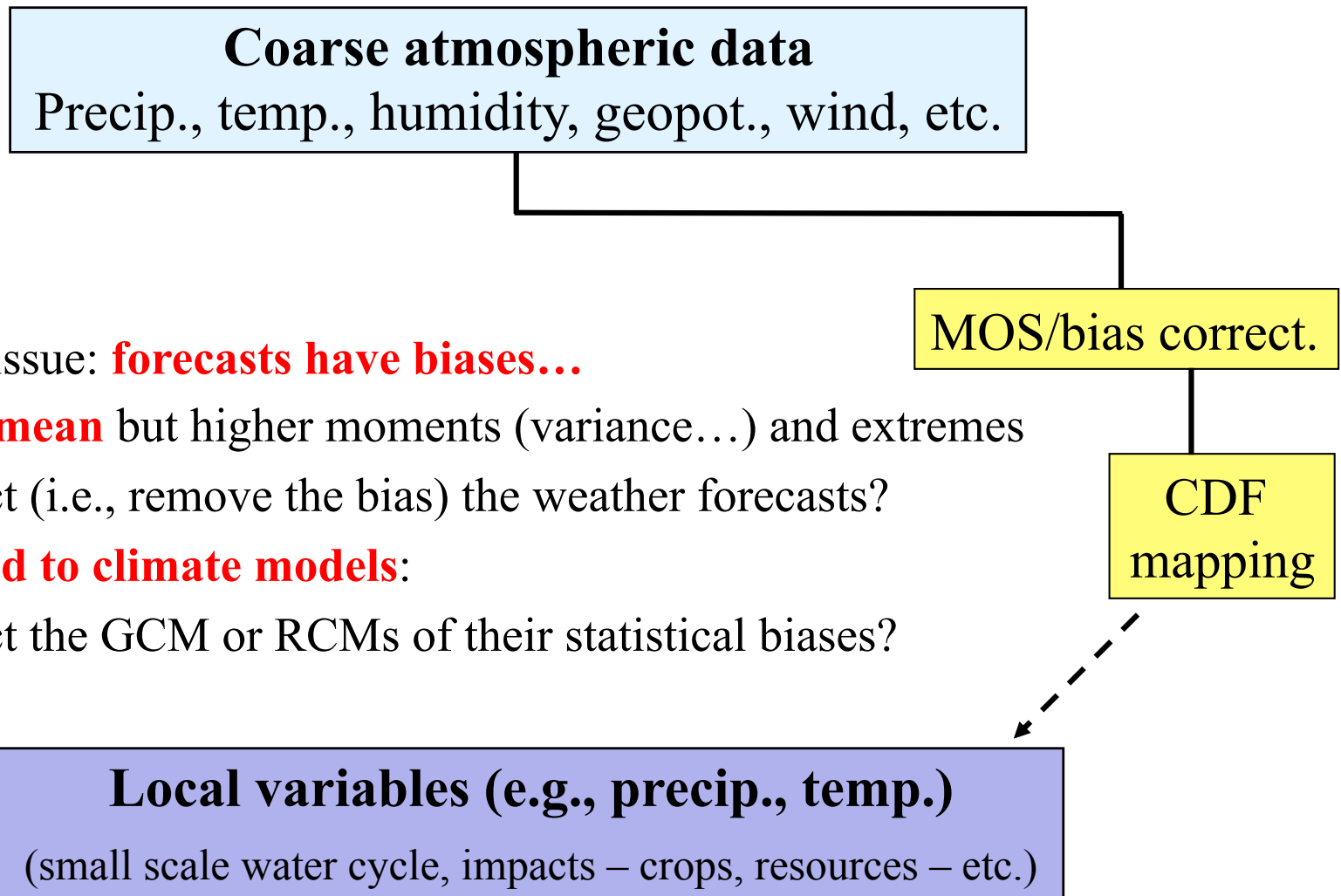
Main statistical downscaling approaches



Main statistical **downscaling** approaches



Main statistical **downscaling** approaches



- Initially, a NWP issue: **forecasts have biases...**

- **Not only the mean** but higher moments (variance...) and extremes

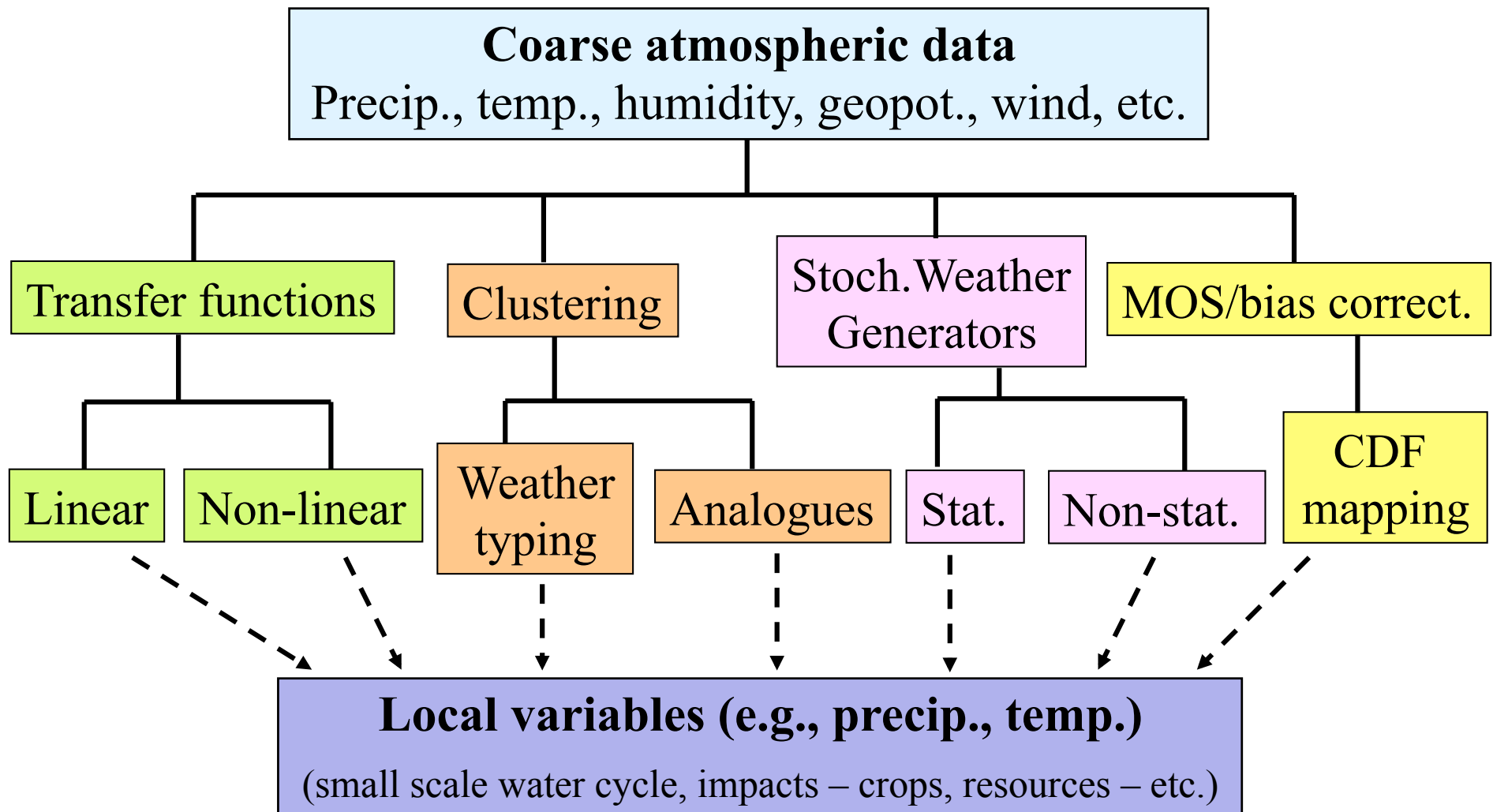
- How to correct (i.e., remove the bias) the weather forecasts?

- Now, **extended to climate models:**

How to correct the GCM or RCMs of their statistical biases?

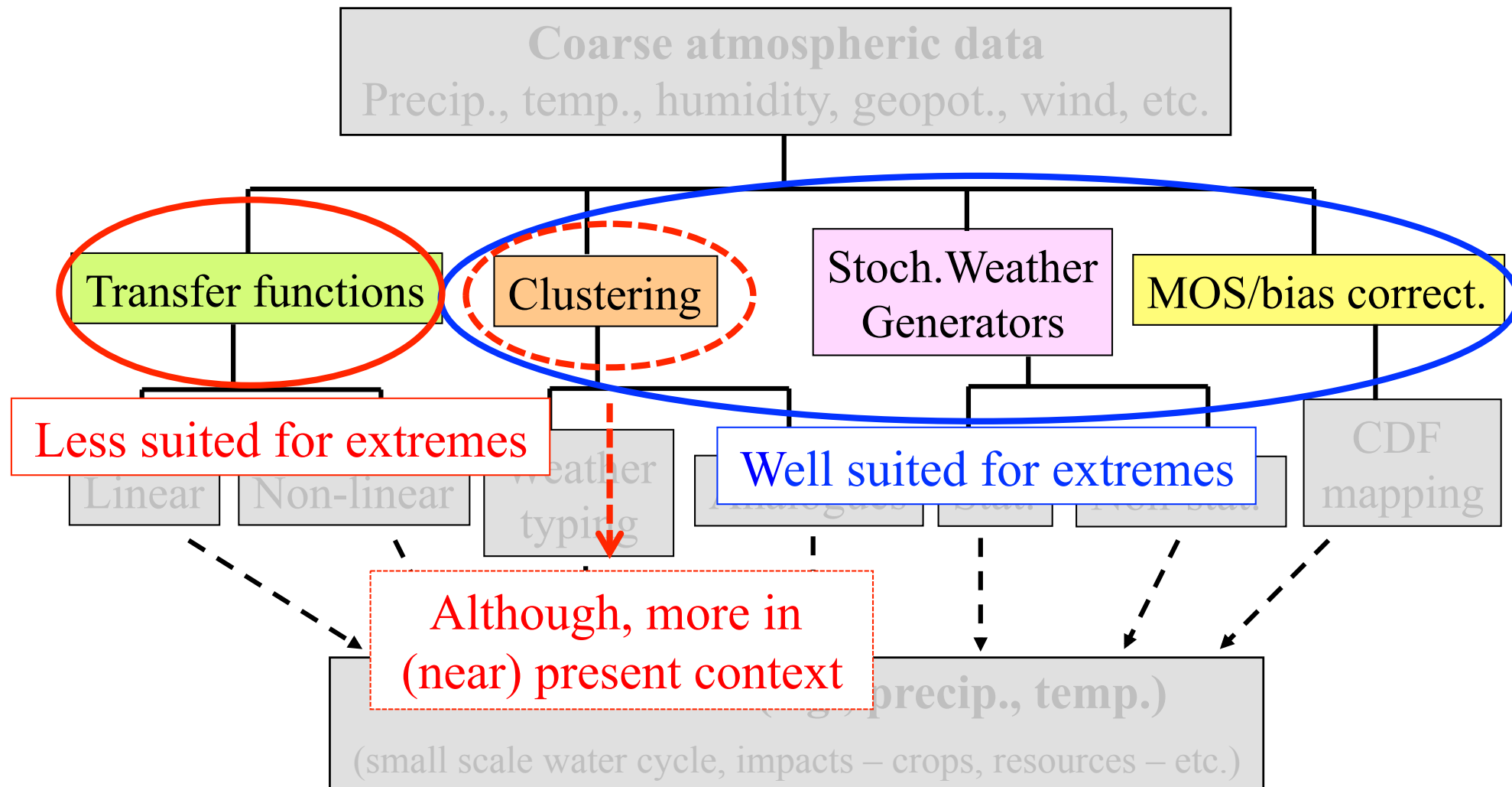
Main statistical **downscaling** approaches

Could also be RCM simulations...



Main statistical **downscaling** approaches

Could also be RCM simulations...





Stochastic weather generators

One illustration with 2 models

VGLM & NN-CMM

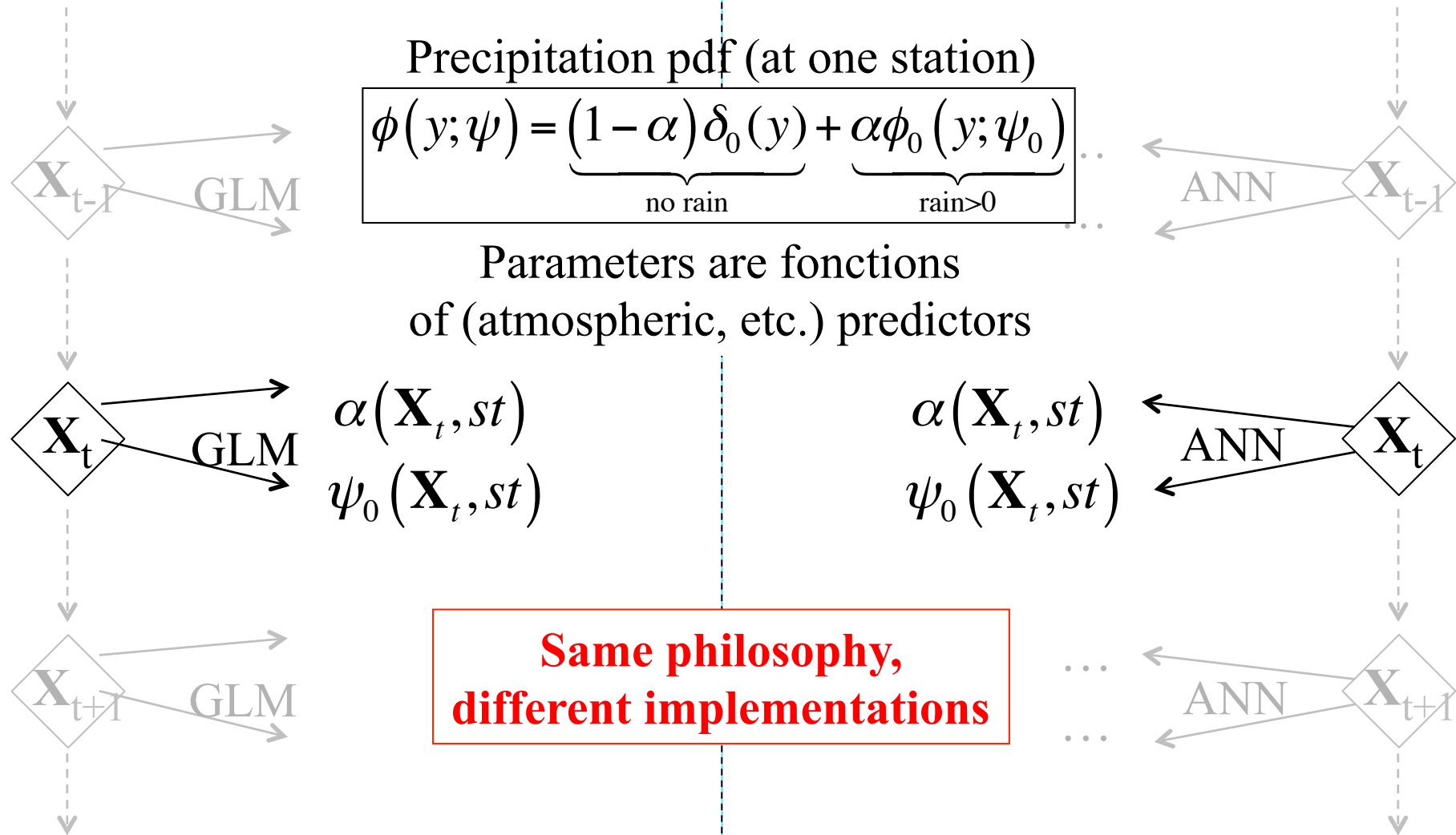
VGLM

&

NN-CMM

Vector Generalized Linear Model

Neural Network – Conditional Mixture Model



Vrac et al. (2007, WRR)
Eden et al. (2014, JGR)

Carreau & Vrac (2011, WRR)

The modelling part of **VGLM**

- Precipitation probability density function (N stations):

$$\begin{aligned}\phi_{\mathbf{Y}_t|\mathbf{X}_t}(\mathbf{y}) &= \prod_{i=1}^N \left[\phi(y_i; \psi_i(\mathbf{X}_t)) \right] \\ &= \prod_{i=1}^N \left[(1 - \alpha_i(\mathbf{X}_t)) \delta_0(y_i) + (\alpha_i(\mathbf{X}_t) \phi_0(y_i; \psi_{0,i}(\mathbf{X}_t))) \right]\end{aligned}$$

with $\alpha_i(\mathbf{X}_t) = \text{Logistic regression}(\mathbf{X}_t)$

$$= \frac{\exp(\mathbf{X}_t' \boldsymbol{\lambda}_i)}{1 + \exp(\mathbf{X}_t' \boldsymbol{\lambda}_i)}$$

and $\phi_0 = \text{Gamma pdf with parameters}$

$$\psi_{0,i}(\mathbf{X}_t) = \begin{cases} k_i(\mathbf{X}_t) = a_0 + a_1 X_1 + \dots + a_p X_p = a_0 + \mathbf{A}\mathbf{X}_t \\ \beta_i(\mathbf{X}_t) = b_0 + b_1 X_1 + \dots + b_p X_p = b_0 + \mathbf{B}\mathbf{X}_t \end{cases}$$

The modelling part of **NN-CMM**

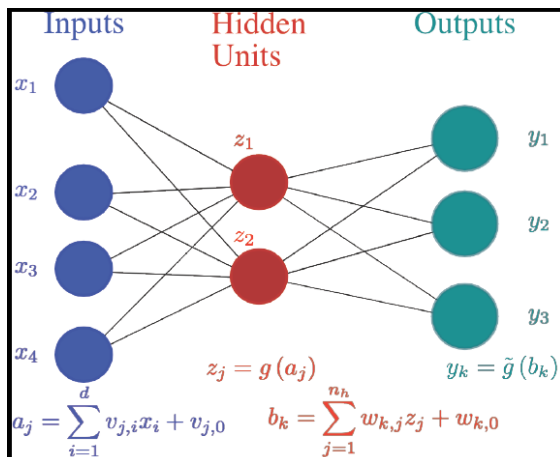
- Precipitation probability density function (N Stations):

$$\phi_{\mathbf{Y}_t|\mathbf{X}_t}(\mathbf{y}) = \prod_{i=1}^N \left[\phi(y_i; \psi_i(\mathbf{X}_t)) \right]$$

$$= \prod_{i=1}^N \left[(1 - \alpha_i(\mathbf{X}_t)) \delta_0(y_i) + \left(\alpha_i(\mathbf{X}_t) \phi_0(y_i; \psi_{0,i}(\mathbf{X}_t)) \right) \right]$$

with $\phi_0(y; \psi_{0,i}(\mathbf{X}_t)) \Leftarrow \sum_{j=1}^m \pi_{i,j}(\mathbf{X}_t) f(y; \theta_{i,j}(\mathbf{X}_t))$

$$\psi_i(\mathbf{x}) = \left(\alpha_i(\mathbf{x}), \left(\pi_{i,j}(\mathbf{x}) \right)_{j=1,\dots,m}, \left(\theta_{i,j}(\mathbf{x}) \right)_{j=1,\dots,m} \right)$$



$$f = \left\{ \begin{array}{l} \text{➤ Gaussian} \\ \text{or} \\ \text{➤ **Log-Normal**} \\ \text{or} \\ \text{➤ Hybrid Pareto} \end{array} \right.$$

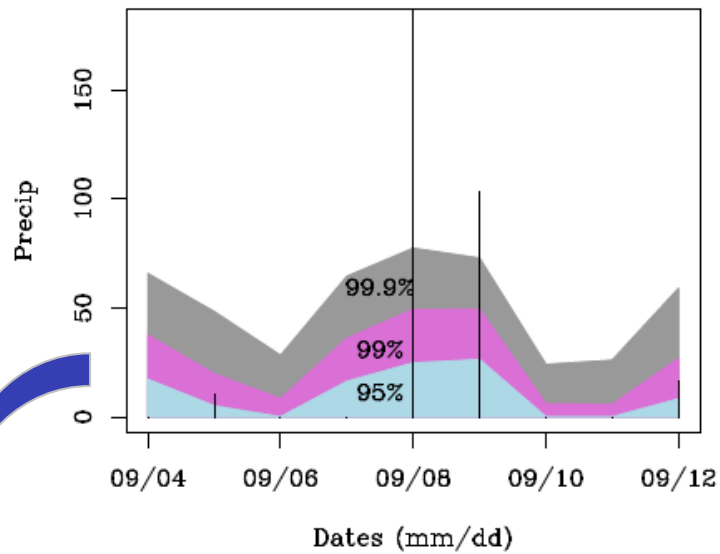
✓ Carreau & Vrac (2011)

✓ Carreau & Bengio (2009a,b)

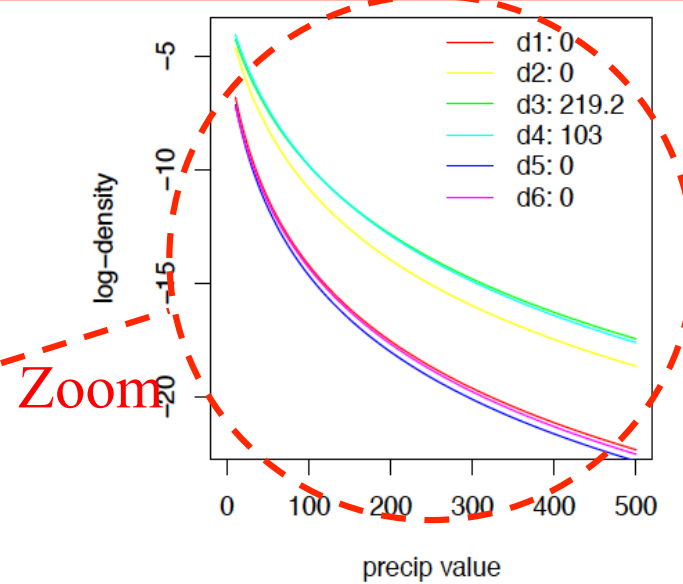
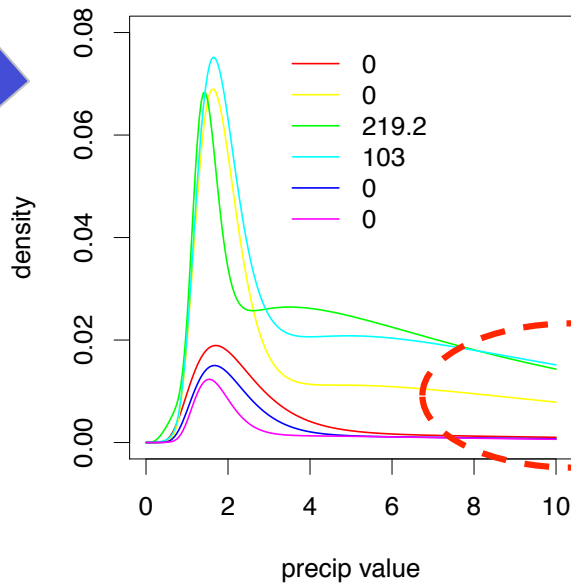
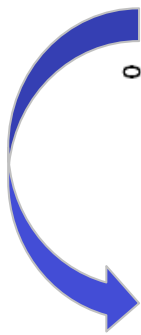
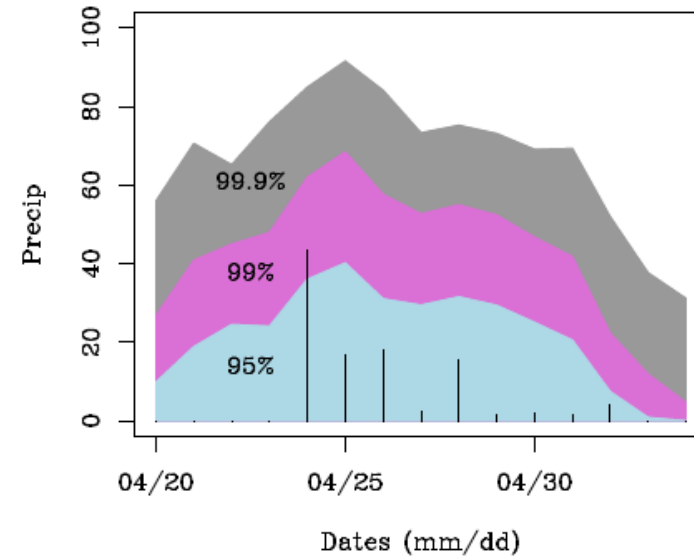
One illustration: Daily pdfs with **NN-CMM-2L**

from Carreau and Vrac (2011)

Spell with the **highest** cum. vol. of rain



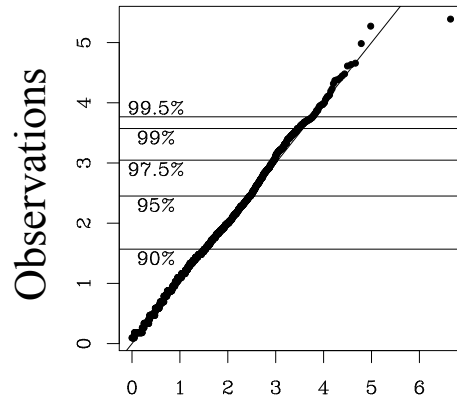
Longest wet spell



NN-CMM-2L vs. (NN-Cond'l) Gamma

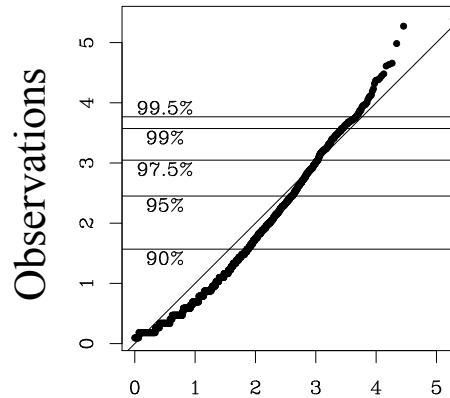
Illustration on the Orange station

QQ-plot (log) **CMM-2L**



Simulations

QQ-plot (log) **Ber-Gamma**

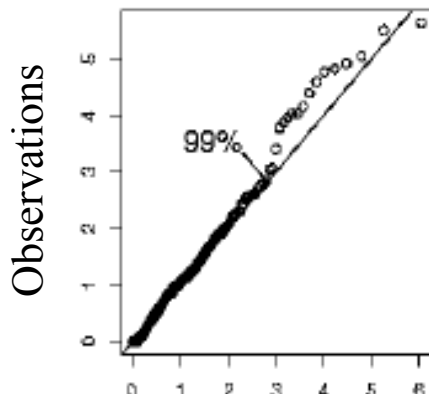


Simulations

Williams (1998)

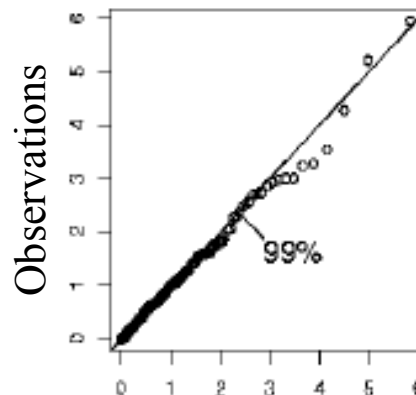
VGLM w/ Gamma: good but... not always enough!

QQplot Charleston



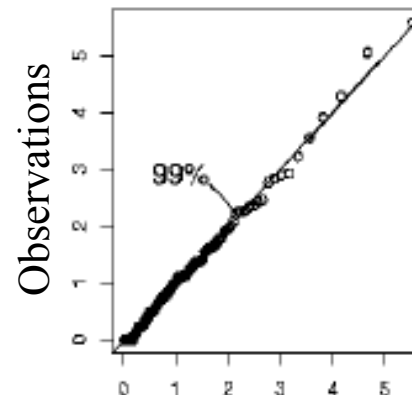
Simulations

QQplot Galva



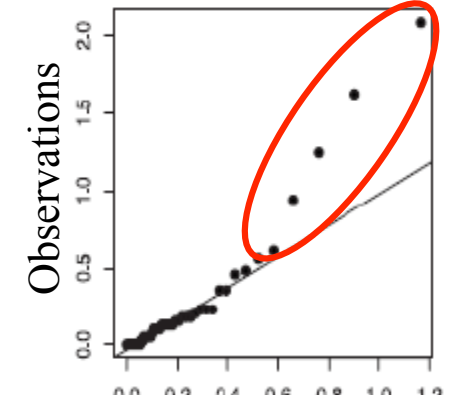
Simulations

QQplot Walnut



Simulations

QQplot Quincy



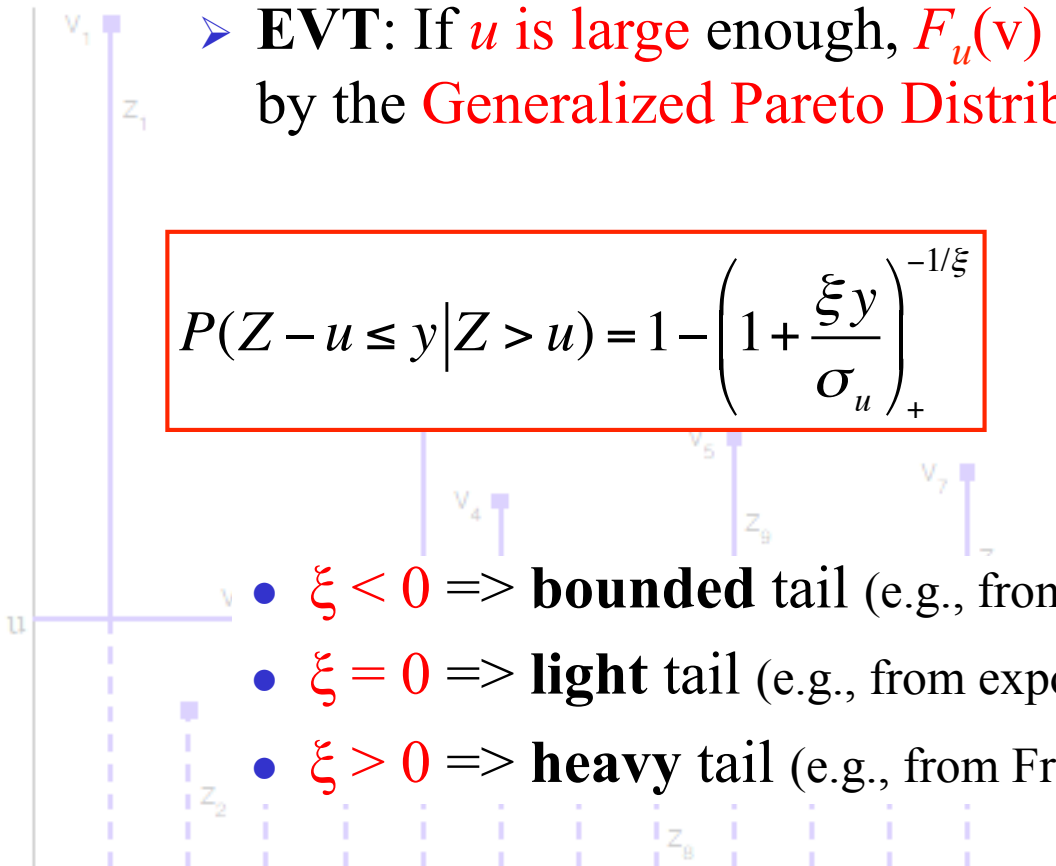
Simulations

Peaks over threshold (POT): Generalized Pareto Distribution (GPD)

- Not simply values higher than the threshold but **excesses**
 - Excess V of the variable Z above threshold u is defined as $Z-u$, given that $Z > u$: $V = Z - u \mid Z > u$
 - **EVT**: If u is large enough, $F_u(v)$ can be approximated by the **Generalized Pareto Distribution (GPD)**

$$P(Z - u \leq y \mid Z > u) = 1 - \left(1 + \frac{\xi y}{\sigma_u} \right)_+^{-1/\xi}$$

- ✓ u = selected threshold
- ✓ σ_u = scale parameter (>0)
- ✓ ξ = shape parameter



- $\xi < 0 \Rightarrow$ **bounded** tail (e.g., from uniform, Weibull, Beta)
- $\xi = 0 \Rightarrow$ **light** tail (e.g., from exponential, Gaussian, Gumbel)
- $\xi > 0 \Rightarrow$ **heavy** tail (e.g., from Fréchet, Student t, Cauchy)



Modelling the **whole** precipitation distribution

Vrac & Naveau (2007)

Stochastic downscaling of precipitation: From dry events to heavy rainfalls, Water Resources Research, 43, W07402, doi: 10.1029/2006WR005308

Wong, Maraun, Vrac, Widmann, Eden, Kent (2014)

Stochastic model output statistics for bias correcting and downscaling precipitation including extremes, Journal of Climate, 27, 6940–6959, doi: <http://dx.doi.org/10.1175/JCLI-D-13-00604.1>

Merging classical and EV distributions in VGLM

- Based on Frigessi et al. (2002):

“*Dynamic mixture model for unsupervised tail estimation without threshold*”

$$\phi_0(y|\psi_0) = c_{\psi_0} \left[\underbrace{(1 - w(y|m, \tau))}_{\text{functional weight}} \underbrace{\Gamma(y|\gamma, \lambda)}_{\text{Gamma pdf}} + w(y|m, \tau) \underbrace{GPD(y|\xi, \sigma, u=0)}_{\text{Generalized Pareto Distribution (GPD) pdf}} \right]$$

with

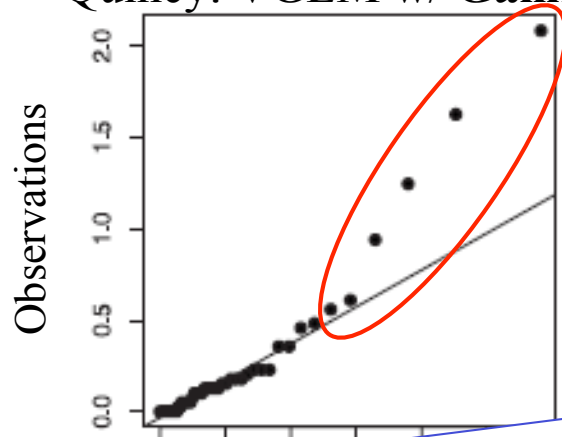
$$w(y|m, \tau) = \frac{1}{2} + \frac{1}{\pi} \arctan\left(\frac{y - m}{\tau}\right)$$

Value where transition from Γ to GPD

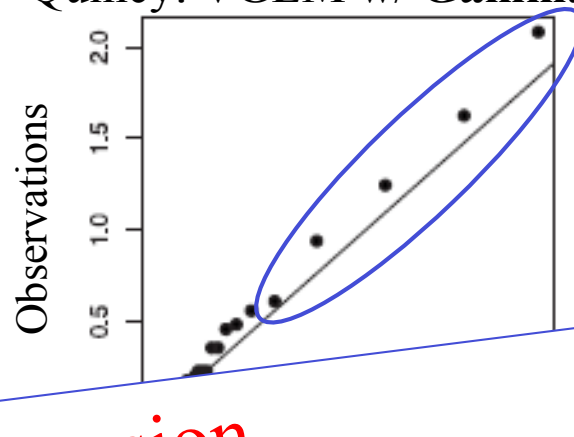
Transition rate

Illustration on two stations

Quincy: VGLM w/ **Gamma**

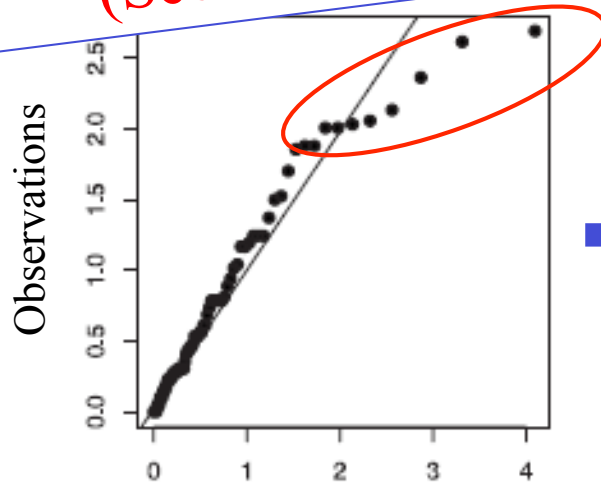


Quincy: VGLM w/ **Gamma&GPD**

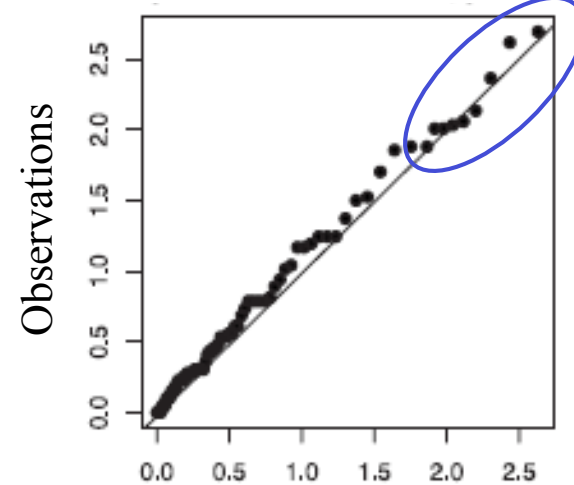


Bivariate extension
(See Vrac, Naveau, Drobinski, 2007, NPG)

Aledo: VGLM w/ **Gamma**



Aledo: VGLM w/ **Gamma&GPD**

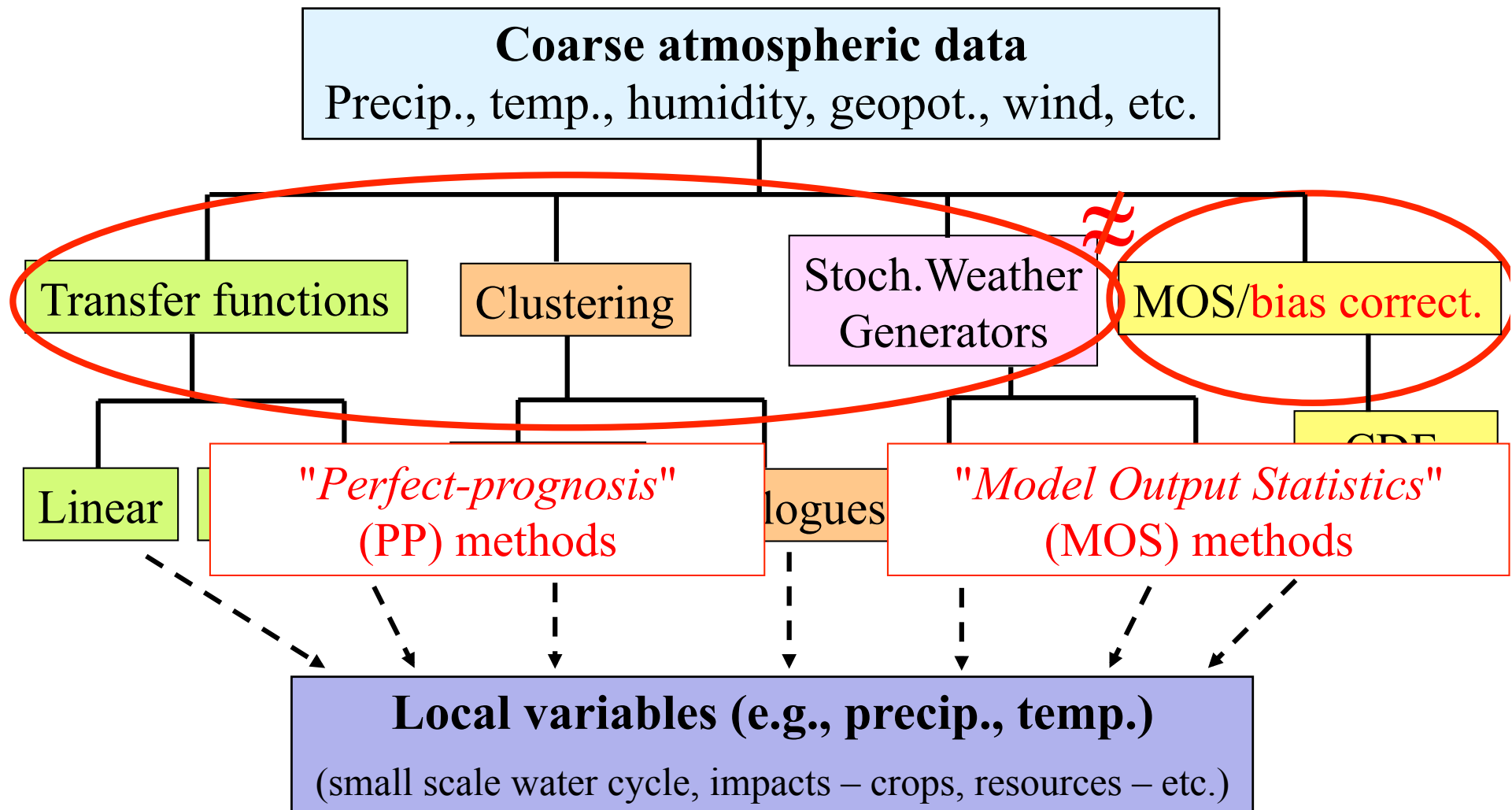


Simulations

Simulations

Main statistical **downscaling** approaches

Could also be RCM simulations...

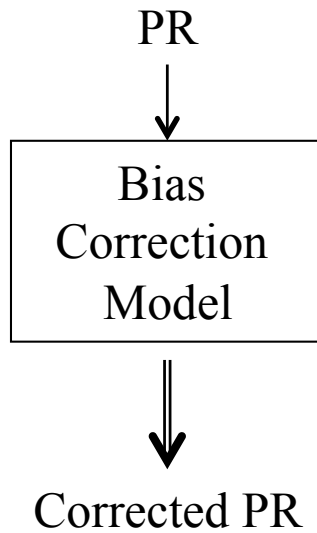


Bias correction



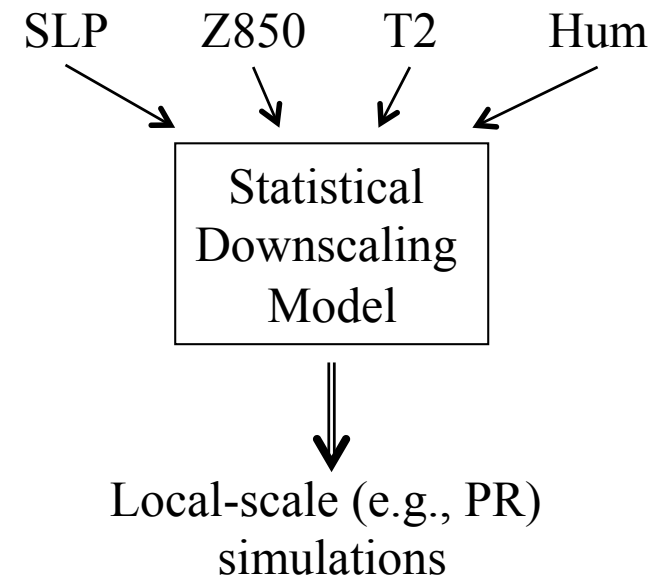
Statistical downscaling

1. ONE predictor



vs.

Several predictors



Bias correction



Statistical downscaling

1. ONE predictor

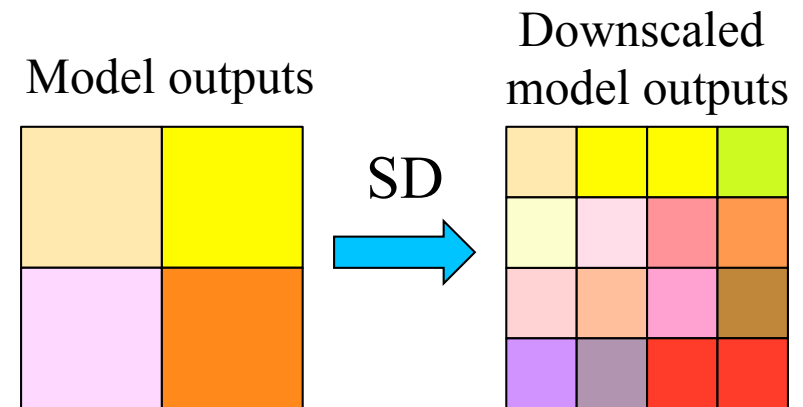
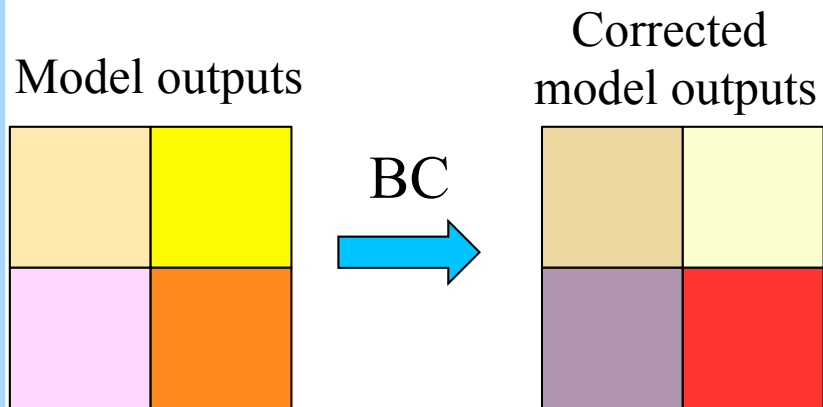
vs.

Several predictors

2. Not necessarily local scale

vs.

Local scale



Bias correction



Statistical downscaling

1. ONE predictor vs.
2. Not necessarily local scale vs.
3. “Model Output Statistics” vs.

Several predictors

Local scale

“Perfect prognosis”

Directly calibrated to link
model outputs & observations

Assumes “perfect” predictors

- calibration **needs temporal matching** between (large-scale) predictors and (local-scale) observations
- Projections based on predictors from GCMs

Bias correction

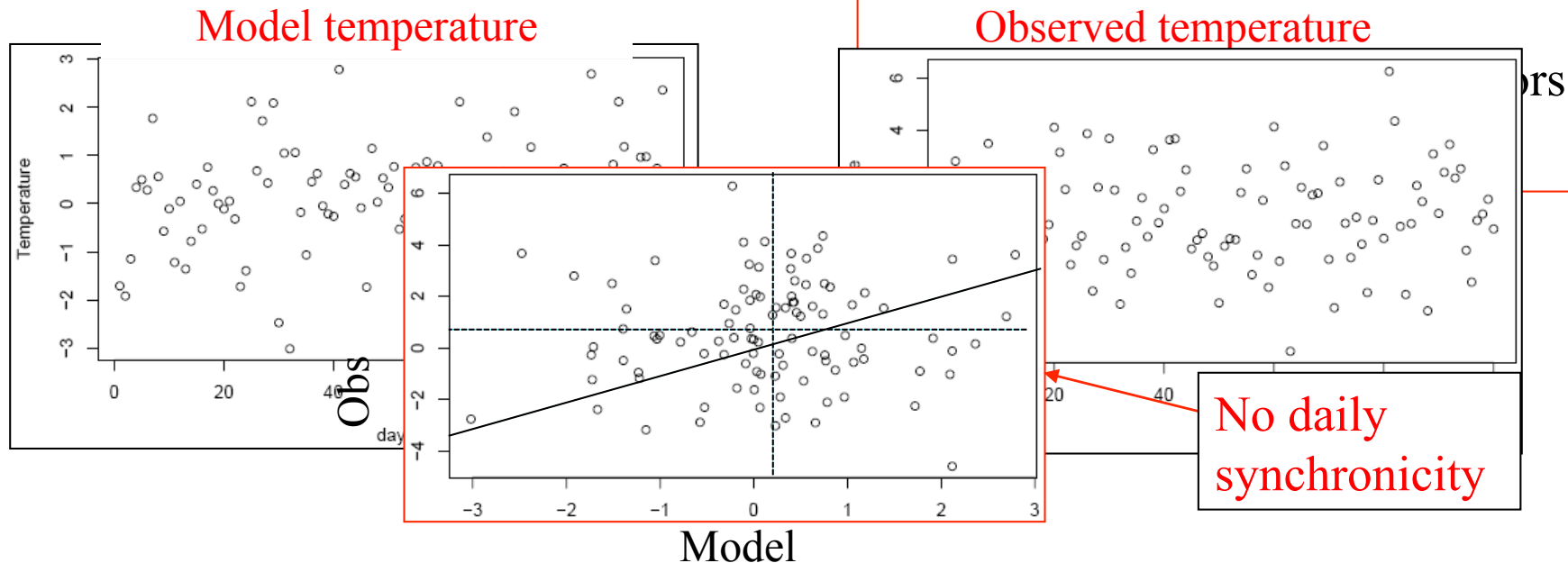


Statistical downscaling

- | | | |
|--------------------------------|-----|---------------------|
| 1. ONE predictor | vs. | Several predictors |
| 2. Not necessarily local scale | vs. | Local scale |
| 3. “Model Output Statistics” | vs. | “Perfect prognosis” |

Directly calibrated to link model outputs & observations

Assumes “perfect” predictors
- calibration **needs temporal matching** between (large-scale) predictors and (local-scale)



Bias correction



Statistical downscaling

1. ONE predictor vs.
2. Not necessarily local scale vs.
3. “Model Output Statistics” vs.

Several predictors

Local scale

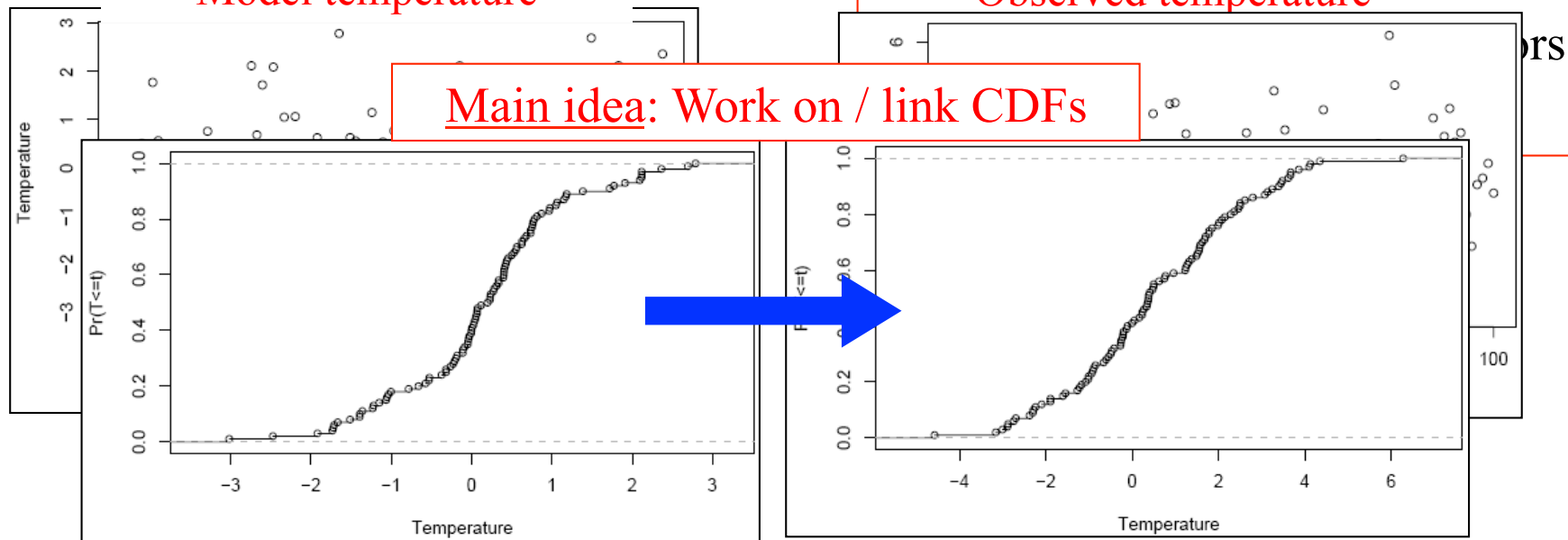
“Perfect prognosis”

Directly calibrated to link model outputs & observations
⇒ Link between the distributions
⇒ No need of temporal matching

Assumes “perfect” predictors
- calibration **needs temporal matching** between (large-scale) predictors and (local-scale)

Model temperature

Observed temperature



Bias correction



Statistical downscaling

1. ONE predictor vs.
2. Not necessarily local scale vs.
3. “Model Output Statistics” vs.
4. **Model-dependant** vs.

- Several predictors
- Local scale
- “Perfect prognosis”
- Same for any model**

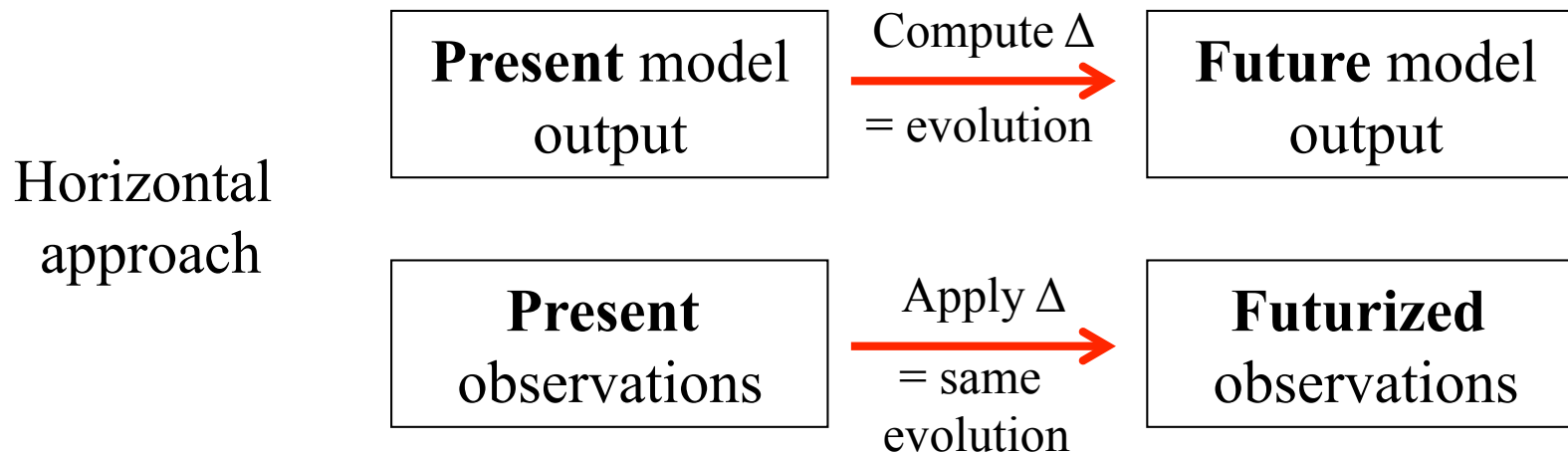


One climate model
⇒ One BC model
(N climate models ⇒ N BC models)

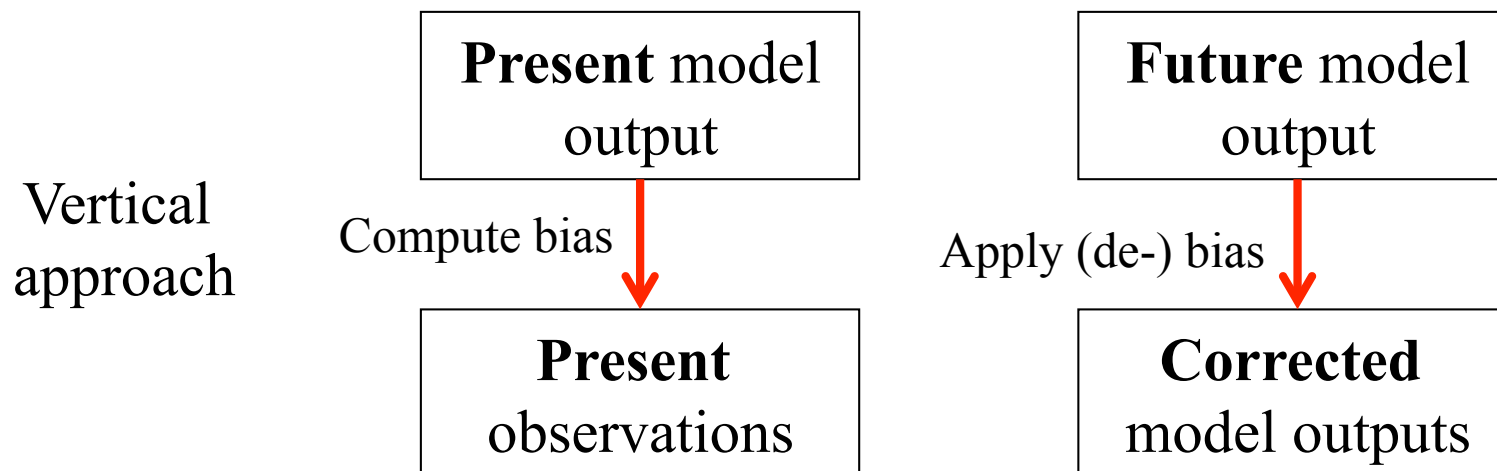
N climate models
⇒ One SD model (calibrated
once, e.g., on reanalyses)

Bias correction: main methods

➤ “Delta”-like methods:



➤ “Quantile-quantile”- or “anomaly”-like methods:



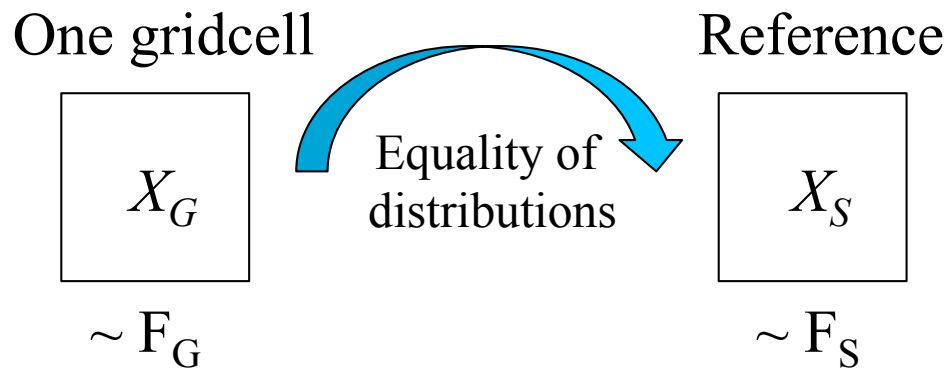
“Quantile-quantile”-like methods

- Classical approach: Quantile-mapping

- Gridcell G; $X_G \sim F_G$; Station S; $X_S \sim F_S$

$$F_S(x_S) = F_G(x_G) \iff x_S = F_S^{-1}F_G(x_G)$$

You want this (red text) points to x_S in the first equation.
You know this (blue text) points to x_G in the first equation.
You obtain this (red text) points to x_S in the second equation.
You know this (blue text) points to x_G in the second equation.



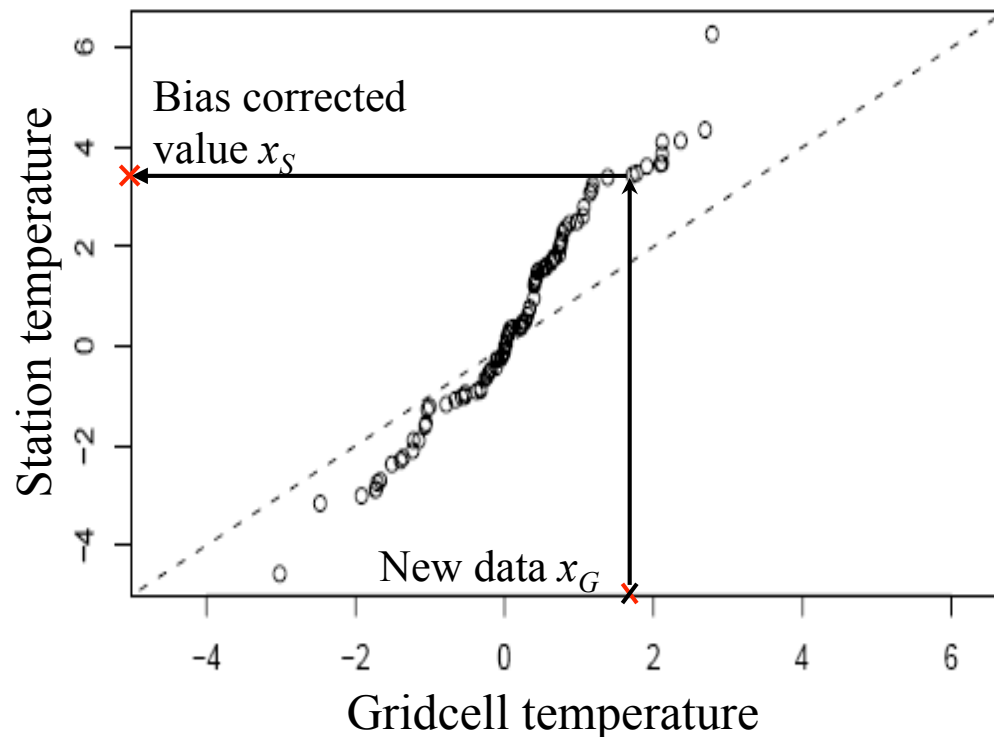
“Quantile-quantile”-like methods

- Classical approach: Quantile-mapping

- Gridcell G; $X_G \sim F_G$; Station S; $X_S \sim F_S$

$$F_S(x_S) = F_G(x_G) \Leftrightarrow x_S = F_S^{-1}F_G(x_G)$$

- Visual interpretation: QQplot (between F_S and F_G)



First paper(s):

Panofsky and Brier (1958);

Haddad and Rosenfeld (1997)

Many variants:

Wood et al. (2004);

Déqué (2007);

Shabalova et al. (2003);

Piani et al. (2010);

etc.

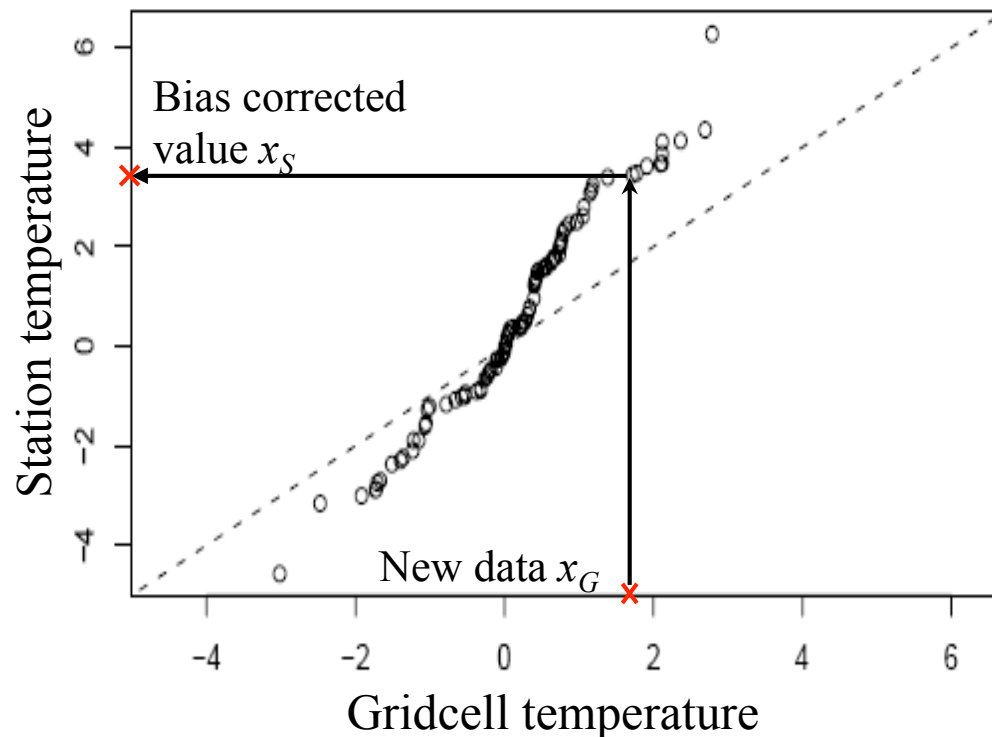
“Quantile-quantile”-like methods

- Classical approach: Quantile-mapping

- Gridcell G; $X_G \sim F_G$; Station S; $X_S \sim F_S$

$$F_S(x_S) = F_G(x_G) \Leftrightarrow x_S = F_S^{-1}F_G(x_G)$$

- Visual interpretation: QQplot (between F_S and F_G)



- 2 main limitations

- Limited to min and max

? What if x_G out of the calib. range ??

- Implicit assumption:

CDF(proj. period) = CDF(calib. period)

? What if the CDF changes ??

“Quantile-quantile”-like methods

Cumulative Distribution Function - transform (CDF-t)

$$x_{Sf} = F_{Sp}^{-1} F_{Gp}(x_{Gf})$$

Classical QQ:

x_{Gf} projected on F_{Gp}

	Present	Future
GCM (1 gridcell)	F_{Gp}	F_{Gf}
Station	F_{Sp}	F_{Sf}

$$X_{Gf} = \{x_{Gf, i}\}_{i=1, \dots, N}$$

QQ from CDF-t:

x_{Gf} projected on F_{Gf}

$$x_{Sf} = F_{Sf}^{-1} F_{Gf}(x_{Gf})$$

“Quantile-quantile”-like methods

Cumulative Distribution Function - transform (CDF-t)

- Formulation (from Vrac et al., 2012):
 - Based on mathematical transformation T applied to LS CDF
 - F_{Sp} Verifies eq.(1) by definition
 - **Assumption: Eq. (2) remains valid in the future**

	Present	Future
GCM (1 gridcell)	F_{Gp} ↓ T	F_{Gf} ⋮ T
Station	F_{Sp}	F_{Sf}

$$T(F_{Gp}(x)) = F_{Sp}(x) \quad (1)$$

Let $x = F_{Gp}^{-1}(u)$ with $u \in [0,1]$

$$\Rightarrow T(u) = F_{Sp}(F_{Gp}^{-1}(u)) \quad (2)$$

$$F_{Sf}(x) = T(F_{Gf}(x)) \Leftrightarrow F_{Sf}(x) = F_{Sp}(F_{Gp}^{-1}(F_{Gf}(x))) \quad (3)$$

“Quantile-quantile”-like methods

Cumulative Distribution Function - transform (CDF-t)

- Formulation (from Vrac et al., 2012):
 - Based on mathematical transformation T applied to LS CDF
 - F_{Sp} Verifies eq.(1) by definition
 - Assumption: F_{Gp} is the large-scale CDF

Advantage:
Changes of the large-scale CDF directly accounted for

GCM (1 gridcell)	F_{Gp}	F_{Gf}
	↓ T	↓ T
Station	F_{Sp}	F_{Sf}

$$T(F_{Gp}(x)) = F_{Sp}(x) \quad (1)$$

Let $x = F_{Gp}^{-1}(u)$ with $u \in [0,1]$

$$\Rightarrow T(u) = F_{Sp}(F_{Gp}^{-1}(u)) \quad (2)$$

$$F_{Sf}(x) = T(F_{Gf}(x)) \Leftrightarrow F_{Sf}(x) = F_{Sp}(F_{Gp}^{-1}(F_{Gf}(x))) \quad (3)$$

✓ **Methodology & evaluations:**

Michelangeli et al. (2009, wind)

Kallache et al. (2011, extreme PR)

Vrac et al. (2012, T&PR)

Vrac & Vaittinada Ayar (2016, combined BC/DS)

Vrac et al. (2016, stochastic/PR, SSR)

Volosciuk et al. (2017, combined BC/DS), ... etc.

✓ **Applications:**

Oettli et al. (2011, T/PR/Rad/evt for crop model)

Colette et al. (2012, BC/RCM)

Tisseuil et al. (2012, river flows)

Vautard et al. (2012, DRIAS)

Vigaud et al. (2013, PR for Indian water resources)

Defrance et al. (2017, Africa PR), ... etc.

✓ **Intercomparison exercises:**

Vaittinada Ayar et al. (2015, EURO/MED-CORDEX)

Gutierrez et al. (2018, VALUE)

Hertig et al. (2018, extremes), ... etc.

Cumulative Distribution Function - transform

Two extensions : **Extremes** & covariates

$$\text{CDF-t} \Rightarrow F_{Sf}(x) = F_{Sp} \left(F_{Gp}^{-1} \left(F_{Gf}(x) \right) \right)$$

➤ **XCDF-t: F's are GPD's** (Kallache, Vrac, Michelangeli, Naveau, 2011, JGR)

$$F_{Sf}(y - c_{fac}) = 1 - \left(1 + \frac{\xi_{Sp}}{\xi_{Gp}} \frac{\sigma_{Gp}}{\sigma_{Sp}} \left[\left(1 + \frac{\xi_{Gf}}{\sigma_{Gf}} y \right)^{\xi_{Gp}/\xi_{Gf}} - 1 \right] \right)^{-(1/\xi_{Sp})}$$

More complex than a GPD...

But... If we assume $\xi_{Gf} = \xi_{Gp}$, then

$$F_{Sf}(y) = 1 - \left(1 + \frac{\xi_{Sp}}{\sigma_{Gf}} \frac{\sigma_{Gp}}{\sigma_{Sp}} (y + c_{fac}) \right)^{-(1/\xi_{Sp})}$$

$\Rightarrow F_{Sf}$ is a **GPD** with $\xi_{Sf} = \xi_{Sp}$ and $\sigma_{Sf} = \sigma_{Gf} \left(\sigma_{Sp} / \sigma_{Gp} \right)$

Cumulative Distribution Function - transform

Two extensions : Extremes & **covariates**

$$\text{CDF-t} \Rightarrow F_{Sf}(x) = F_{Sp} \left(F_{Gp}^{-1} \left(F_{Gf}(x) \right) \right)$$

- **XCDF-t: F's are GPD's** (Kallache, Vrac, Michelangeli, Naveau, 2011, JGR)

If we assume $\xi_{Gf} = \xi_{Gp}$, then

$$\Rightarrow F_{Sf} \text{ is a GPD with } \xi_{Sf} = \xi_{Sp} \text{ and } \sigma_{Sf} = \sigma_{Gf} \left(\sigma_{Sp} / \sigma_{Gp} \right)$$

- **Inclusion of covariate information:**

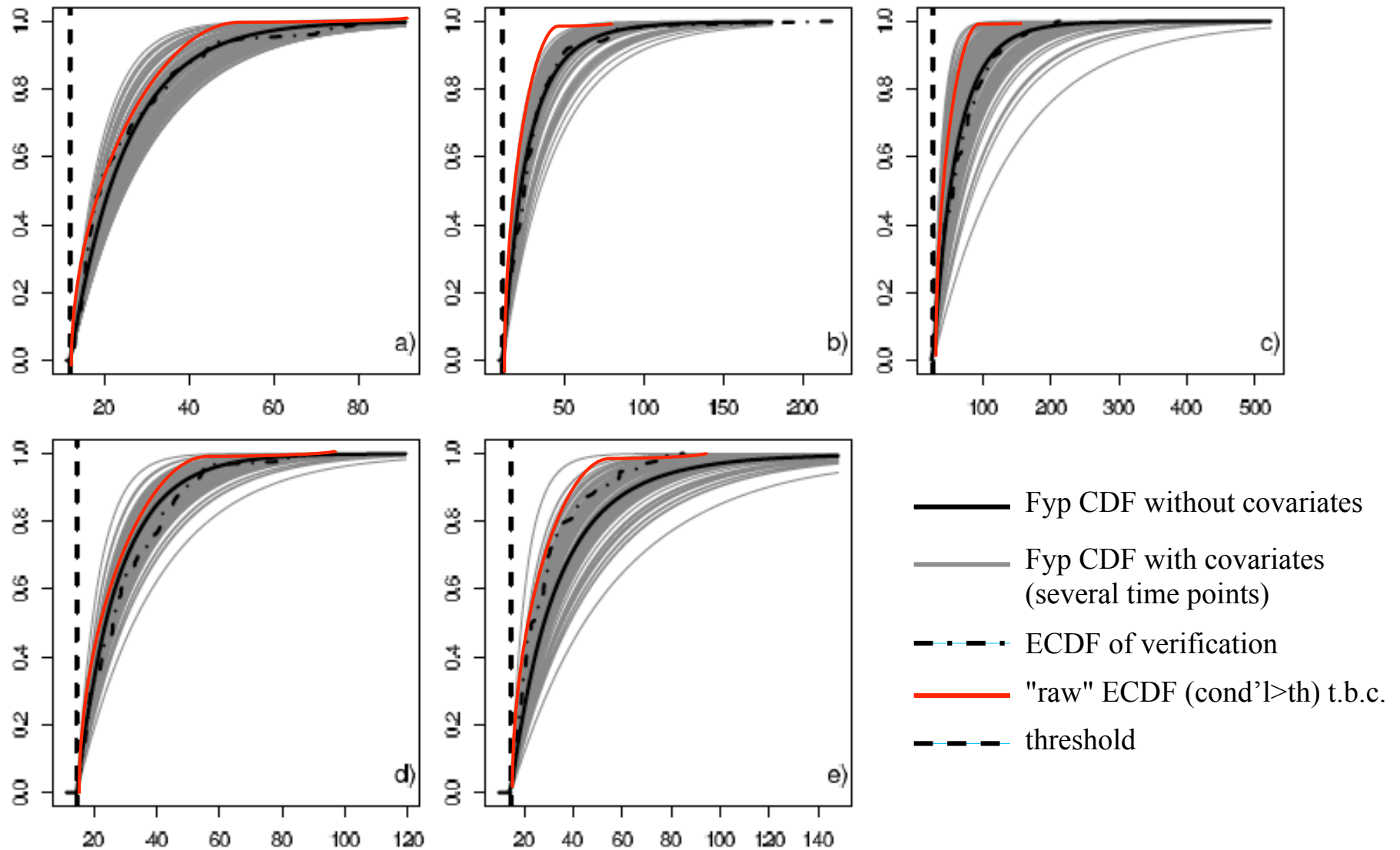
$$F(.) = GPD(\sigma, \xi)$$

$$\sigma_t = \exp \left(a_0 + a_1 \text{cov}_1^t + \dots + a_n \text{cov}_n^t \right)$$

- ✓ The parameters may have different link functions
- ✓ The covariates vary with time (e.g., as in Clim.Ch.)
 \Rightarrow one CDF $F_{Sf}^t(.)$ per time step !
- ✓ Feature similar to (cond'l) SWG

eXtreme CDF-t (XCDF-t)



One illustration



XCDF-t on precipitation at 5 French stations.

Calib=1951-1985, Proj= 1986-1999, x-axis in mm/d, y-axis in probability

Conclusions (some) on downscaling & BC

- **Many** (and many) **models and applications** of downscaling & BC
 - My favorite ones:
 - ✓ *Stochastic WGs*: cond'l event-wise variability/uncertainty
 - ✓ *MOS / Bias correction*: DS of CDFs from CDFs
 - Choice of the predictors is a major issue in Stat. DS
 - Non-stationarity ( the SWGs should not explode )
 - Applying Stochastic WGs to GCMs *may be* better than to RCMs
- **RCMs vs. SDMs**: Not a conflict => complementary approaches
 - ⇒ Both have pros & cons
- **There is not one good SDM for all variables and regions**
 - ⇒ Different skills according to regions/variables/applications, etc.
 - ⇒ **Use ensembles** if possible!

Commercial break (well, it's free)

- Some **R packages** developed for Stochastic downscaling & BC:

- **NHMixt** (Vrac & Naveau, 2007, Wong et al., 2014)
 - ✓ Statistical mixture model Gamma & GPD
 - ✓ Inclusion of covariates
- **condmixt** (Carreau & Vrac et al., 2011)
 - ✓ ANN-Conditional mixture model
 - ✓ Various distributions (Gaussian, Log-N, hybrid Pareto)
- **McSIM** (Bechler, Vrac, Bel, 2015)
 - ✓ Spatial models for extreme (maxima) downscaling
 - ✓ Max-stable processes
- **CDVineCopulaConditional** (Bevacqua et al., 2017)
 - ✓ Copula-based model for compound events
 - ✓ Multivariate dependence
- **CDFt** (Vrac et al., 2012) & **XCDFt** (Kallache et al., 2011) → **Run at IPSL**
 - ✓ Bias correction
 - ✓ Also for excesses (i.e., GPD)
- **EC-BC** (Vrac and Friederichs, 2015) → **R²D²** (Vrac, 2018)
 - ✓ For multivariate properties of BC/DS
 - ✓ Post-processing of any 1d method



<http://www.r-project.org>
Or my website

Some perspectives: dependences

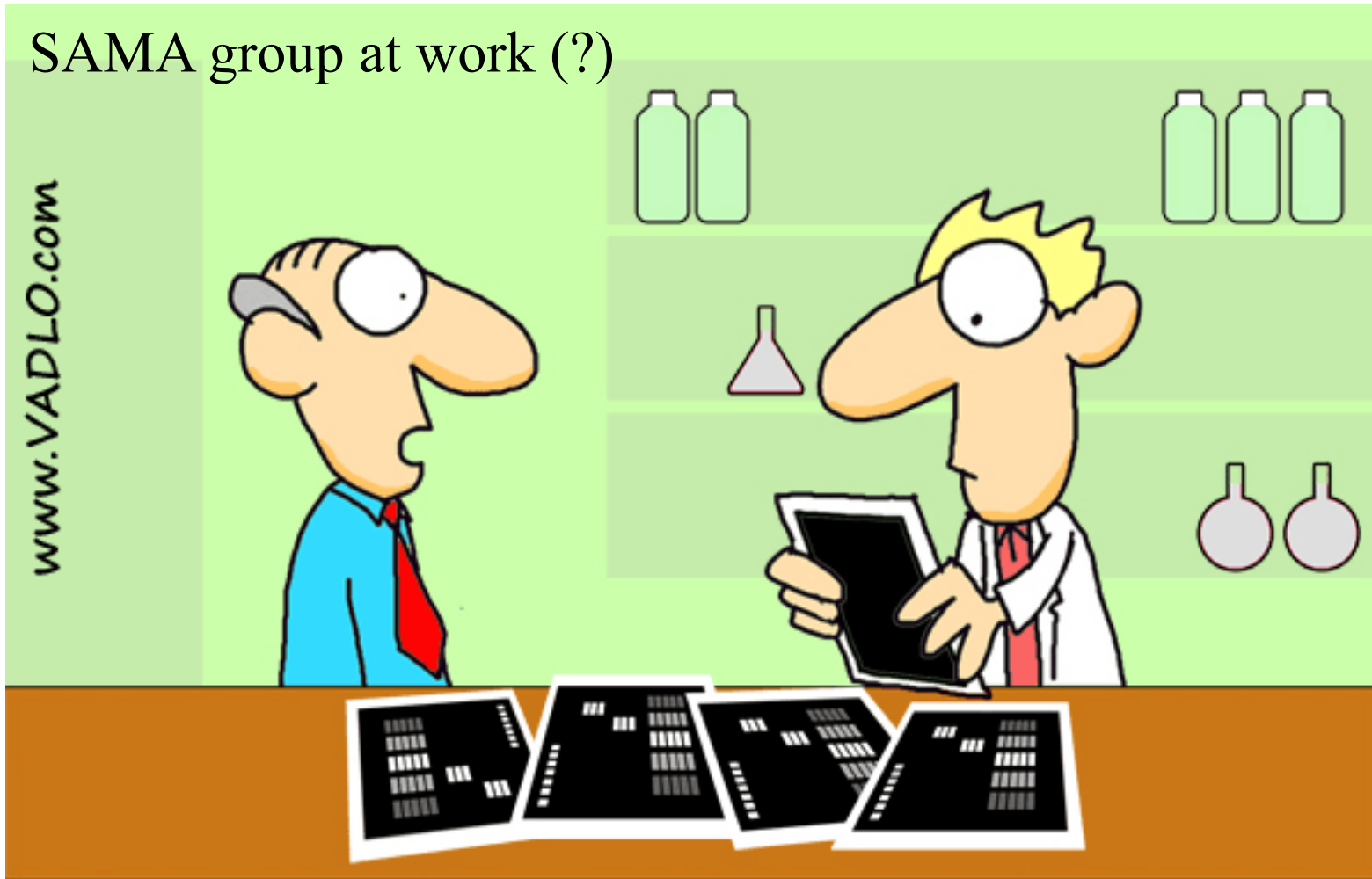
- **Spatial/multi-sites** dependence structure, e.g.:
 - Spatial VGLM (Chandler and Wheater, 2002)
 - ✗ Stationary spatial dependences / Limited number of locations
 - “*Hybrid Spatial Downscaling*” (HSD) for extreme fields (Bechler et al., 2015)
 - ✓ Combines RCM with geostatistical cond'l simulations / DS anywhere in the region
 - “*Rank Resampling for Distributions and Dependences*” (R2D2, Vrac, 2018)
 - ✓ 2-step approach: univariate SDM or BC + dependence reconstruction
- **Multi-variables** dependence structure, e.g.:
 - Copula-based model for compound events (Bevacqua et al., 2017)
 - ✓ Parametric model (=> can be limited in dimensions)
 - R2D2 (Vrac, 2018)
 - ✓ Non-parametric approach / Multi-sites & multi-variables / even in high-dimension
- Multidimensional (sites and/or variables) **dependence of extremes**
 - More complex / Non-parametric approach (e.g., Naveau et al., 2014)
 - Needed to improve/extend "Extreme Event Attribution", impact studies, etc.

Some perspectives: dependences

- Spatial/multi-sites dependence structure, e.g.:
 - Spatial VGLM (Chandler and Wheater, 2002)
 - ✗ Stationary spatial dependences / Limited number of locations
 - “Hybrid Spatial Downscaling” (HSD) for extreme fields (Bechler et al., 2015)
 - ✓ Combines RCM with geostatistical cond’l simulations / DS anywhere in the region
 - “Peak Resampling for Distributions and Dependences” (R2D2, Vrac, 2018)
 - ✓
 - The questions are not necessarily technical anymore:
 - What do you **trust** (or not) in the model outputs?
 - What do you want to **correct/preserve**?
- Mult
 - C
 - ✓ Parametric model (=> can be limited in dimensions)
 - R2D2 (Vrac, 2018)
 - ✓ Non-parametric approach / Multi-sites & multi-variables / even in high-dimension
- Multidimensional (sites and/or variables) **dependence of extremes**
 - More complex / Non-parametric approach (e.g., Naveau et al., 2014)
 - Needed to improve/extend "Extreme Event Attribution", impact studies, etc.

Thank you...

SAMA group at work (?)



“Data don’t make any sense,
we will have to resort to statistics.”