# Improving cross-domain brain tissue segmentation in fetal MRI with synthetic data[⋆]

Vladyslav Zalevskyi[1,2], Thomas Sanchez[1,2], Margaux Roulet[1,2], Jordina Aviles Verddera[3], Jana Hutter[3,4], Hamza Kebiri[1,2], and Meritxell Bach Cuadra[2,1]
vladyslav.zalevskyi@unil.ch

[1] Department of Radiology, Lausanne University Hospital and University of Lausanne (UNIL), Lausanne, Switzerland
[2] CIBM Center for Biomedical Imaging, Switzerland
[3] Department for Early Life Imaging, School of Biomedical Engineering & Imaging Sciences, King's College London, London, UK
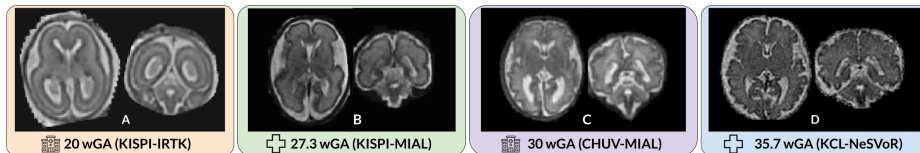[4] Smart Imaging Lab, Diagnostic Radiology, FAU Erlangen-Nuremberg, Erlangen, Germany

**Abstract.** Segmentation of fetal brain tissue from magnetic resonance imaging (MRI) plays a crucial role in the study of *in utero* neurodevelopment. However, automated tools face substantial domain shift challenges as they must be robust to highly heterogeneous clinical data, often limited in numbers and lacking annotations. Indeed, high variability of the fetal brain morphology, MRI acquisition parameters, and super-resolution reconstruction (SR) algorithms adversely affect the model's performance when evaluated out-of-domain. In this work, we introduce FetalSynthSeg, a domain randomization method to segment fetal brain MRI, inspired by SynthSeg. Our results show that models trained solely on synthetic data outperform models trained on real data in out-of-domain settings, validated on a 120-subject cross-domain dataset. Furthermore, we extend our evaluation to 40 subjects acquired using low-field (0.55T) MRI and reconstructed with novel SR models, showcasing robustness across different magnetic field strengths and SR algorithms. Leveraging a generative synthetic approach, we tackle the domain shift problem in fetal brain MRI and offer compelling prospects for applications in fields with limited and highly heterogeneous data.

**Keywords:** Domain shifts · segmentation · fetal brain · MRI · synthetic data · low-field MRI

**Fig. 1.** Domain shifts across data splits in fetal SR MRI. (GA in weeks, site-SR). A & C - pathological, B & D - neurotypical.

## 1 Introduction

Fetal brain magnetic resonance imaging (MRI) is a growing diagnostic tool for studying neurodevelopment in fetuses [1–4]. Despite its potential, creating automated pipelines for fetal MRI faces challenges due to the limited availability of annotated datasets and data heterogeneity. The fetal brain undergoes significant morphological changes during gestation and can be severely altered by pathologies, which complicates its automatic analysis [5–7]. Additionally, datasets encounter distribution shifts from variations in acquisition sites, scanners, and imaging protocols [8, 9]. Utilizing super-resolution-reconstructed (SR) volumes addresses issues like fetal motion artefacts [10] and low through-plane acquisition resolution, but introduces additional heterogeneity in texture, tissue contrast, intensity values and other reconstruction artefacts [11] (see Figure 1). A recent study [12] on the FeTA 2022 MICCAI challenge revealed significant performance drops in white matter (WM), gray matter (GM), and ventricles segmentation when models were tested on diverse clinical datasets, highlighting the impact of domain shifts on the analysis of fetal MRI.

Numerous techniques exist to mitigate domain shifts, including domain adaptation [13, 14], transfer learning [15], meta-learning [16], style transfer [17] and data harmonization [18]. However, these methods typically rely on the availability of at least target domain images which is challenging in domains with limited data. Other approaches that use multi-centre learning [19, 20] require multiple training domains which are costly and labour-intensive to acquire. In fetal brain imaging, where new SR algorithms and MRI scanners introduce significant diversity [21], available domains may not sufficiently capture the required variability for generalizable models. Recent advancements in single-source domain generalization (SSDG) involve techniques such as global intensity non-linear augmentation (GIN) and causal interventions [22]. These modifications eliminate spurious correlations, enhancing model robustness to variations in image intensities and textures. In a related study [23], authors analyzed frequency's effect on domain discrepancy, using a mixed frequency spectrum for self-supervised augmentation. While effective, limitations may arise in inducing spatial and intensity transformations, especially with diverse SR algorithms introducing artefacts associated with skull stripping or the inclusion of extra-cerebral tissue.

A promising alternative approach involves synthetic data generation based on segmentation maps [24, 25], achieving domain generalization through domain

randomization [26], only requiring labels but no images. Leveraging shape information from segmentations, these models introduce diverse spatial and intensity transformations, along with flexible artefact simulations, mitigating many factors causing domain shifts in MRI.

In our study, we delve into exploring how synthetic generative models can be used to construct a diverse fetal brain dataset for training segmentation models. Our contributions are the following: i) We adapt the domain randomization of Billot et al. [24] to fetal brain MRI, accommodating specific fetal anatomical properties, acquisition artefacts and heterogeneity due to fetal brain development and SR algorithms; ii) We show that our method, trained only using synthetic data, performs better than models trained using real data when evaluated out-of-domain and performs on par with state-of-the-art SSDG algorithms; iii) We extend our evaluation to low-field (0.55T) MRI data, showing the robustness of our approach to unseen magnetic field strength and SR algorithms.

## 2   Methodology

### 2.1   Data

Various datasets are used in our experiments to validate the SSDG efficacy of the models we explore. These datasets come from multiple institutions and were acquired using MRI scanners from various manufacturers, with different field strengths, acquisition parameters and reconstructed with different SR algorithms. The acquisition details are given in Table 1.

**FeTA dataset.** We used the publicly available data from the MICCAI 2022 FETA challenge (KISPI) [27, 28]. It consists of 80 subjects, among which are 40 reconstructed using MIALSRTK [29] and 40 using Simple-IRTK [30]. Expert annotators delineated seven tissue labels (external cerebrospinal fluid (eCSF), GM, WM, ventricles, cerebellum, deep GM, and brainstem). The ethical committee of the Canton of Zurich, Switzerland approved the prospective and retrospective studies that collected and analysed KISPI MRI data (Decision numbers: 2017-00885, 2016-01019, 2017-00167).

**Clinical 1.5T dataset.** Additionally, we include a proprietary clinical dataset, named `CHUV`, containing 40 subjects reconstructed with MIALSRTK [29] and manually annotated following the FeTA protocol [27]. Data was retrospectively collected from acquisitions done between January 2013 to April 2021. All images were anonymized. This dataset is part of a larger research protocol approved by the ethics committee of the Canton de Vaud (decision number CER-VD 2021-00124) for re-use of their data for research purposes and approval for the release of an anonymous dataset for non-medical reproducible research and open science purposes. This private dataset was used to evaluate methods submitted for the FETA 2022 challenge and a detailed description of it can be found in [12].

**Clinical 0.55T dataset.** We also evaluate our model on 40 neurotypical cases acquired on a low-field scanner from another centre, referred to as `KCL`. Subjects were SR reconstructed twice, using two novel methods, NeSVoR [11] and SVRTK

[31]. There are no available manual segmentations for this dataset. Fetal MRI was acquired at Kings College London and approved for sharing with interested academic researchers around the world by the Ethics Committee London Bromley (Ethics code 21/LO/0742). The data has been acquired during a prospective single-center study and has been fully anonymised in line with local procedures

## 2.2   FetalSynthSeg

Our generative model is inspired by SynthSeg [24] which is based on domain randomization [26]. SynthSeg [24] leverages image segmentation as a structural prior, integrating randomization across resolution, intensity, contrast, and spatial distortions. This approach yields a diverse dataset that is well-suited for training models capable of robustly handling these sources of variation. We adapt this method to fetal brain segmentation by introducing some crucial changes related to the tissue generation classes (see Figure 2).
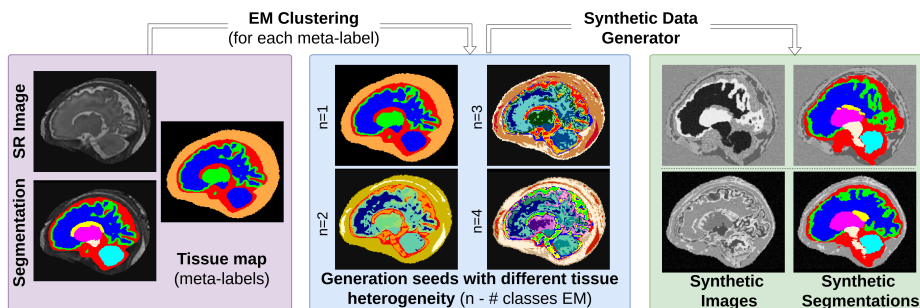
**Table 1.** Dataset properties.

| Site | Scanner | Acquisition Parameters | SR | Res. $(mm^3)$ | GA $(weeks)$ | $N_n$ | $N_p$ |
|------|---------|------------------------|-----|---------------|--------------|-------|-------|
| KISPI* | GE Signa Discovery MR450/MR750 (1.5T/3T) | ssFSE TR: 2500-3500/120 ms 0.5 x 0.5 x 3.5 $mm^3$ | mial | $0.5^3$ | 20-34 | 25 | 15 |
| | | | irtk | $0.5^3$ | 20-35 | 24 | 16 |
| CHUV* | Siemens MAGNETOM Aera (1.5T) | HASTE TR/TE: 1200/90 ms 1.1 x 1.1 x 3 $mm^3$ | mial | $1.1^3$ | 21-35 | 25 | 15 |
| KCL | Siemens MAGNETOM FREE.MAX (0.55T) | HASTE TR/TE: 2500/106 ms 1.5 x 1.5 x 4.5 $mm^3$ | svrtk | $0.8^3$ | 21-35 | 40 | 0 |
| | | | nesvor | | | | |

*FeTA Data [27]                              $N_n$: neurotypical, $N_p$: pathological

Instead of directly using target segmentations as generation classes, our approach first introduces an intermediate seed generation step. In it, we initiate synthetic image generation by defining four primary *meta-labels*: CSF, WM, GM, and skull with surrounding tissue. Then we employ the expectation–maximisation (EM) algorithm [32] for intensity clustering within each meta-class, resulting in 1-4 subclasses per meta-class. In a second step, these subclasses serve as inputs for synthetic data generators, producing images and segmentations that faithfully reflect the observed heterogeneity in SR and those intrinsic to fetal brains, for example, related to WM maturation [33]. By leveraging tissue-type subclasses rather than original segmentation labels during synthetic image generation, we mitigate reliance on artificial contrast disparities between original labels, a critical consideration given the heterogeneity within a single class. This strategy is similar to the one used in [24] for cross-domain cardiac MRI

**Fig. 2.** Synthetic image generation framework. Original segmentation labels are merged to create a 4-meta-label tissue map (CSF, WM, GM, skull). EM clustering then divides each meta-label into 1 to 4 subclasses, capturing tissue heterogeneity. A generative model uses these split meta-labels to produce synthetic images.

segmentation, which showed the importance of such splitting in representing the target structures with different levels of heterogeneity. Following subclass creation, voxel intensities are independently sampled from Gaussian distributions with randomly selected means and standard deviations. Random artefact corruptions, including bias field simulation, Gaussian blur and noise addition, are applied to introduce common noise and artefacts prevalent in SR images. The generative model applies a battery of spatial transformations (affine and elastic) to simulate spatial distortions. Table S1 in the Supplementary Material provides all detailed parameters of the generative model. We use offline image generation, creating 200 synthetic images per real image to train the model based on purely synthetic images, which is referred to as `FetalSynthSeg` (or `_synth` as a suffix) further in the text. Code and pre-trained models will be released upon acceptance of the paper.

### 2.3   Segmentation model

**Architecture.** Based on Billot et al. [24] and Valabregue et al. [34], our study employed a 3D U-Net with five levels, featuring instance normalization, max-pooling, and upsampling operations in the expanding path. Each level includes a $3 \times 3 \times 3$ kernel convolutional layer with LeakyReLU activation, except the final layer using softmax activation. Starting with 32 feature maps, the initial layer doubles and halves after max-pooling and upsampling layers, respectively. Skip connections facilitate information flow between the contracting and expanding paths. We use the same architecture and training hyperparameters across all datasets and splits to ensure comparability.

**Pre-processing.** During model training, both real and synthetic images undergo identical pre-processing steps before being fed to the model, including resampling to $0.5 \times 0.5 \times 0.5$ mm$^3$ with centre crop and crop-padding to 256x256x256 (when needed), random contrast adjustment via gamma transformation, random affine transformations (scaling, rotation, shearing, and translation), random

Gaussian noise and smoothing with detailed hyperparameters presented in Supplementary Material Table S1. Subsequently, image intensities are normalized between 0 and 1 via min-max normalization.

**Training.** Models are trained using Adam (LR $= 10^{-3}$) on a combination of Dice and Cross-Entropy losses [34] with a `ReduceLROnPlateau` scheduler [35] (factor $= 0.1$, patience $= 10$) for up to 500 epochs (batch size $= 1$). Training halts on persistent validation dice plateau in the last 10 epochs. Internal validation used 5 randomly selected cases per split, ensuring 35 real and 7,000 synthetic images ($200 \times 35$) for the training of baseline and FetalSynthSeg respectively.

### 2.4   Experimental settings

In our experiments, we compare `FetalSynthSeg` model to i) a baseline model trained on real images and labels (denoted from now on as `baseline`), as well as to ii) `fit_nnUnet`[5], FeTA 2022 challenge winner, an ensemble of nnUnet models trained with all 80 images from KISPI using the GIN SSDG approach [22]. We aim to show that our model using only synthetically generated images can outperform models trained on real images when tackling out-of-domain generalization and reach comparable performance to SSDG SotA domain generalization approaches.

**Experiment 1 - High-field generalization.** We perform the first domain generalization experiment by using the FeTA data as well as the clinical 1.5T dataset. This setting replicates the evaluations carried out in the FeTA challenge 2022 [12]. We consider three data splits, defined by images acquired at specific sites and reconstructed with specific SR, namely `KISPI - mial`, `KISPI - IRTK` and `CHUV - mial` (each containing 35 real images for the baseline training, and 7000 synthetic images for synth training). Two models are trained on each split, one using the original FeTA data and labels, and another using the synthetic data generated from the labels only. This yields a total of six models (denoted as `<Site>_<SR>_base/synth` in figures). The models trained on one (`Site`, `SR`) data split are tested on the two other (`Site`, `SR`) pairs. The mean dice score (`mdsc`) and 95th percentile Hausdorff distance of each model variant are reported (95th HD in Supplementary Material), and models are compared using the non-parametric Wilcoxon rank-sum test with Bonferroni correction. The p-value for statistical significance is set to 0.05.
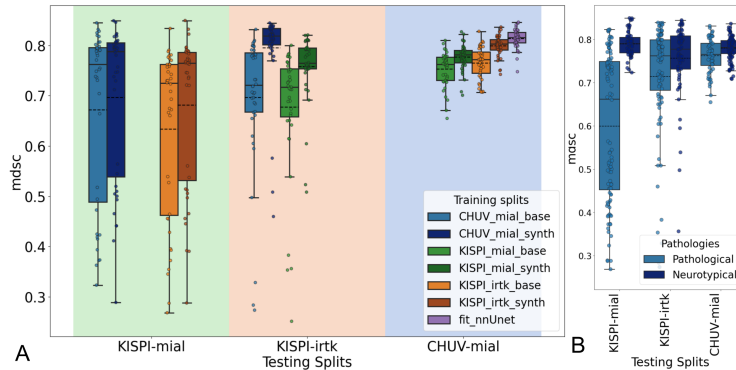
**Experiment 2 - Low-field generalization.** We assess our trained model's adaptability to new, unseen data with a dataset comprising 40 neurotypical subjects acquired at `KCL` on a low-field (0.55T) MRI scanner and reconstructed either using SVRTK [30] or NeSVoR  [11] (none of them used in FeTA data). We select the top-performing model variant trained on original data and the counterpart trained on synthetic data, both derived from the same data split. As no ground truth is available, we evaluate model predictions by comparing tissue volume growth through gestational age (GA) with FeTA reference data. A second-order polynomial fit with confidence interval is evaluated.

---

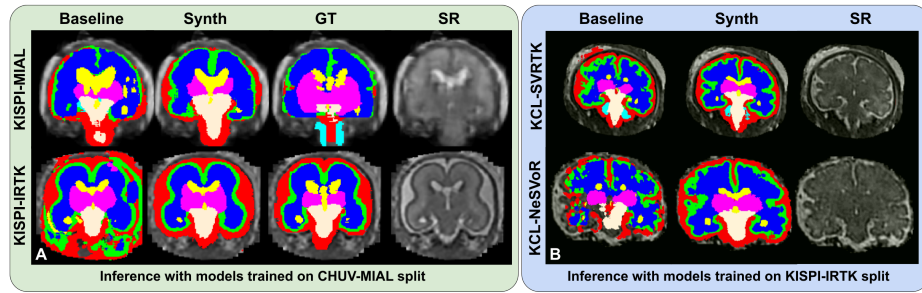[5] The model is available on the FeTA challenge DockerHub page [12].

## 3   Results

**High-field generalization.** A comparison of the out-of-distribution model predictions is shown in Figure 3. Models trained on synthetic data consistently outperformed those trained on original images across all out-of-domain testing splits. Statistical tests confirmed the significance of the differences in mean Dice scores, with p-values $< 0.00005$ for all paired comparisons between corresponding `baseline` and `FetalSynthSeg` models trained and tested on the same split. The same trend is observed with the 95th percentile Hausdorff distances, available in Supplementary Figure S1. Additionally, we compared the segmentation accuracy of our model with the FeTA 2022 challenge winners on the `CHUV-mial` dataset, which was not used for `fit_nnUnet` training. We highlight that we reached a close albeit slightly lower performance to their solution, although our models were trained on half the amount of data (as they rely on only 35 subjects).
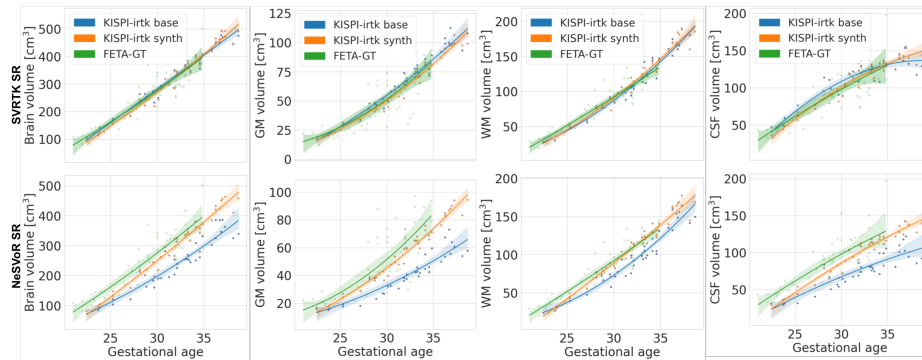
The lowest Dice scores were obtained for pathological cases, particularly evident in the `KISPI-mial` split (see Figure 3B). This discrepancy can be attributed to the fact that approximately 24% of the `KISPI-mial` dataset exhibits poor SR quality often in the severe pathological cases, as noted in [27] and illustrated in Figure 4 (top row). Nonetheless, we noted qualitative improvements in the `FetalSynthSeg` model compared to the `baseline` model, as illustrated in Figure 4 (and in more detail in Figure S2 in the Supplement). Differences in skull stripping by IRTK and MIAL algorithms lead to erroneous segmentation of the skull and surrounding tissue using `baseline` models. However, the synthetic model is more robust to SR-induced domain shifts and artefacts and avoids these errors, even though it was trained on the segmentations from the same split.



**Fig. 3.** Comparison of the out-of-distribution performance of the segmentation models (mdsc - mean Dice score across all tissues). **(A)** `baseline` (light) vs `FetalSynthSeg` (dark). Data split: `KISPI-mial` (green), `KISPI-irtk` (red), `CHUV-mial` (blue). See Figure S1 from the appendix for a comparison with in-distribution performance as well as results split by gestational age. The dashed horizontal line inside the boxplot corresponds to the mean value. **(B)** Pathological vs neurotypical, aggregated across all models.

**Fig. 4.** Cross-domain model inference qualitative results. **(A)** Model trained on `CHUV-mial` and tested on `KISPI-irtk`/`KISPI-mial`. **(B)** Model trained on `KISPI-irtk` and tested on `KCL-svrtk`/`KCL-nesvor`.



**Fig. 5.** Segmented tissue volumes vs GA for `KCL` data reconstructed with SVRTK (top row) and NeSVoR (bottom row). `KISPI_irtk_base` (blue) and `KISPI_irtk_synth` (orange) model predictions are compared to FeTA reference values (green) which are based on the ground truth segmentation of 40 healthy subjects selected across all splits. Lines are second-order polynomial fit and a corresponding shaded area is a confidence interval. See Figure S3 in the Supplement for all tissues evaluation.

**Low-field generalization.** Segmented tissue volumes (total brain, GM, WM and CSF volumes) as a function of GA for `KCL-svrtk` and `KCL-nesvor` segmentations are illustrated in Figure 5. The volume growth curves obtained from SVRTK reconstructions remain within the confidence interval for both `Fetal-SynthSeg` and `baseline` models. However, a notable deviation is observed in all estimated tissue volumes for NeSVoR reconstructions predicted by `KISPI-irtk baseline` model and to a lesser extent from our proposed `FetalSynthSeg` model. This discrepancy highlights a substantial domain gap within SR algo-

rithms, resulting in an underestimation of expected tissue volumes that do not occur on a closer domain to the KISPI-irtk of SVRTK reconstructions. Remarkably, the model trained on synthetic data demonstrates greater robustness to this domain shift compared to the model trained on real data and exhibits minimal deviation from the expected values calculated on the FeTA dataset while qualitatively showing a superior performance as seen in Figure 4B. The baseline model struggles with correct tissue segmentation on NeSVoR reconstructions due to out-of-domain appearance, while performance is improved on SVRTK reconstructions, which exhibit a smaller SR domain gap with IRTK reconstructions. While cortex topology could still be more precise, the synthetic model shows consistent qualitative performance across both scenarios.
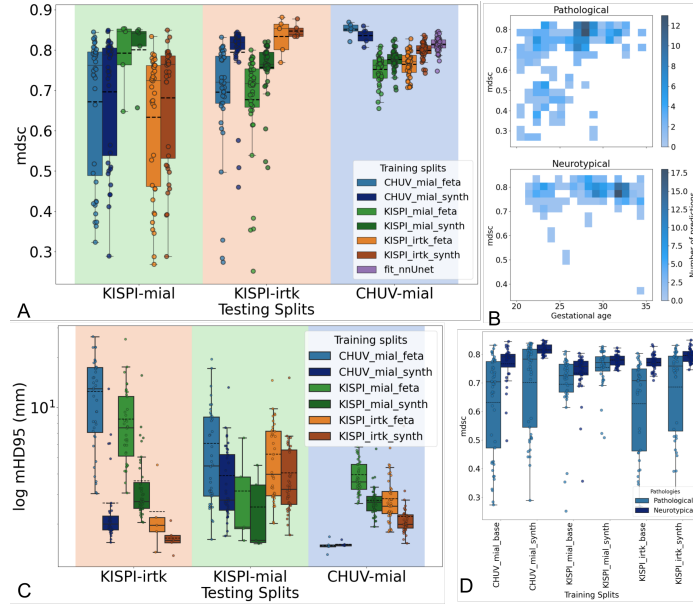
## 4   Conclusion

Our study demonstrates that `FetalSynthSeg` allows robust fetal brain tissue segmentation across datasets with significant domain shifts. We showed how strong randomization of spatial and intensity properties during the synthetic image generation helps models overcome differences caused by MRI acquisition variations and super-resolution reconstruction. Even with half the data, our approach achieved performance close to state-of-the-art SSDG models trained for fetal brain segmentation. Compared to models trained solely on real data, those trained exclusively on synthetic data showed superior performance, particularly in cases of novel SR algorithms or images acquired at different field strengths. The generalization of segmentation models to low-field MRI is of utmost significance, offering an avenue to enhance fetal MRI accessibility in underserved cohorts and low-income regions, by providing a cost-effective diagnostic solution. Our findings suggest that synthetic data can mitigate performance drops caused by limited data and diverse imaging conditions, offering promising applications in fields with highly heterogeneous data, such as fetal imaging.

# Bibliography

[1] O. M. Benkarim *et al.*, "Toward the automatic quantification of in utero brain development in 3d structural MRI: A review," *Human Brain Mapping*, vol. 38, p. 2772–2787, Feb. 2017.

[2] A. Jakab *et al.*, "Emerging magnetic resonance imaging techniques in open spina bifida in utero," *European Radiology Experimental*, vol. 5, no. 1, 2021.

[3] J. Aviles Verdera *et al.*, "Reliability and feasibility of low-field-strength fetal MRI at 0.55 t during pregnancy," *Radiology*, vol. 309, no. 1, 2023.

[4] M. C. Cortes-Albornoz *et al.*, "MR insights into fetal brain development: what is normal and what is not," *Pediatric Radiology*, pp. 1–11, 2024.

[5] J. Wu *et al.*, "Age-specific structural fetal brain atlases construction and cortical development quantification for chinese population," *NeuroImage*, vol. 241, 2021.

[6] C. M. Pfeifer *et al.*, "MRI depiction of fetal brain abnormalities," *Acta Radiologica Open*, vol. 8, Dec. 2019.

[7] F. Vahedifard *et al.*, "Automatic ventriculomegaly detection in fetal brain MRI: A step-by-step deep learning model for novel 2d-3d linear measurements," *Diagnostics*, vol. 13, no. 14, 2023.

[8] H. Guan and M. Liu, "Domain adaptation for medical image analysis: A survey," *IEEE Transactions on Biomedical Engineering*, no. 3, 2022.

[9] R. Lin *et al.*, "Cross-age and cross-site domain shift impacts on deep learning-based white matter fiber estimation in newborn and baby brains," *arXiv preprint arXiv:2312.14773*, 2023.

[10] T. Sanchez *et al.*, "FetMRQC: Automated quality control for fetal brain MRI," in *Perinatal, Preterm and Paediatric Image Analysis*, (Cham), pp. 3–16, Springer Nature Switzerland, 2023.

[11] J. Xu *et al.*, "NeSVoR: Implicit neural representation for slice-to-volume reconstruction in MRI," *IEEE Transactions on Medical Imaging*, 2023.

[12] K. Payette *et al.*, "Multi-center fetal brain tissue annotation (FeTA) challenge 2022 results," *arXiv preprint arXiv:2402.09463*, 2024.

[13] P. de Dumast *et al.*, "Synthetic magnetic resonance images for domain adaptation: Application to fetal brain tissue segmentation," in *2022 IEEE 19th International Symposium on Biomedical Imaging (ISBI)*, IEEE, 2022.

[14] Z. Xu *et al.*, "Asc: Appearance and structure consistency for unsupervised domain adaptation in fetal brain mri segmentation," in *Medical Image Computing and Computer Assisted Intervention – MICCAI*, 2023.

[15] N. Karani *et al.*, "A lifelong learning approach to brain mr segmentation across scanners and protocols," in *MICCAI*, 2018.

[16] D. Li *et al.*, "Learning to generalize: meta-learning for domain generalization," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI Press, 2018.

[17] K. Zhou *et al.*, "Domain generalization with mixstyle," in *International Conference on Learning Representations*, 2021.

[18] F. Hu *et al.*, "Image harmonization: A review of statistical and deep learning methods for removing batch effects and evaluation metrics for effective harmonization," *NeuroImage*, vol. 274, 2023.

[19] D. Li *et al.*, "Episodic training for domain generalization," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, (Los Alamitos, CA, USA), pp. 1446–1455, IEEE Computer Society, 2019.

[20] Q. Liu *et al.*, "Shape-aware meta-learning for generalizing prostate MRI segmentation to unseen domains," in *MICCAI*, 2020.

[21] J. P. Marques *et al.*, "Low-field MRI: An MR physics perspective," *Journal of magnetic resonance imaging*, vol. 49, no. 6, pp. 1528–1542, 2019.

[22] C. Ouyang *et al.*, "Causality-inspired single-source domain generalization for medical image segmentation," *IEEE Transactions on Medical Imaging*, vol. 42, no. 4, p. 1095–1106, 2023.

[23] H. Li *et al.*, "Frequency-mixed single-source domain generalization for medical image segmentation," in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2023*, Springer Nature Switzerland, 2023.

[24] B. Billot *et al.*, "Synthseg: domain randomisation for segmentation of brain scans of any contrast and resolution," *arXiv:2107.09559*, 2021.

[25] B. Billot *et al.*, "Robust machine learning segmentation for large-scale analysis of heterogeneous clinical brain MRI datasets," *Proceedings of the National Academy of Sciences*, vol. 120, no. 9, 2023.

[26] J. Tremblay *et al.*, "Training deep networks with synthetic data: Bridging the reality gap by domain randomization," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018.

[27] K. Payette *et al.*, "An automatic multi-tissue human fetal brain segmentation benchmark using the fetal tissue annotation dataset," *Scientific Data*, vol. 8, no. 1, 2021.

[28] K. Payette and A. Jakab, "Fetal tissue annotation dataset feta," 2021.

[29] S. Tourbier *et al.*, "Medical-image-analysis-laboratory/mialsuperresolutiontoolkit: MIAL super-resolution toolkit v2.0.1." *Zenodo*, 2020.

[30] M. Kuklisova-Murgasova *et al.*, "Reconstruction of fetal brain MRI with intensity matching and complete outlier removal," *Medical image analysis*, vol. 16, no. 8, pp. 1550–1564, 2012.

[31] A. U. Uus *et al.*, "Automated 3d reconstruction of the fetal thorax in the standard atlas space from motion-corrupted MRI stacks for 21–36 weeks ga range," *Medical image analysis*, vol. 80, p. 102484, 2022.

[32] A. P. Dempster *et al.*, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 39, no. 1, pp. 1–22, 1977.

[33] H. Lajous *et al.*, "A fetal brain magnetic resonance acquisition numerical phantom (FaBiAN)," *Scientific Reports*, vol. 12, May 2022.

[34] R. Valabregue *et al.*, "Comprehensive analysis of synthetic learning applied to neonatal brain MRI segmentation," *arXiv:2309.05306*, 2023.

[35] A. Paszke *et al.*, "PyTorch: An imperative style, high-performance deep learning library," 2019.

## Supplementary material



**Fig. S1.** Detailed segmentation evaluation. **(A)** Comparison between models including in-domain evaluation on the 5 validation cases. **(B)** Aggregated performance across all models stratified by GA of subjects. **(C)** Boxplot presenting mean 95th-percentile Hausdorff-Distance scores across all tissues for our experiments on a log scale. **(D)** Per-model mean dice score results of models evaluated on `KISPI-mial` split with the distinction between pathological and neurotypical cases.

**Table S1.** Hyperparameters of the synthetic generator and augmentations. Only parameters deviating from the default ones used by the fully randomized generative model without any tissue priors [24] are reported. Intensities are in $[0, 255]$ range, rotations in degrees and spatial parameters in mm.

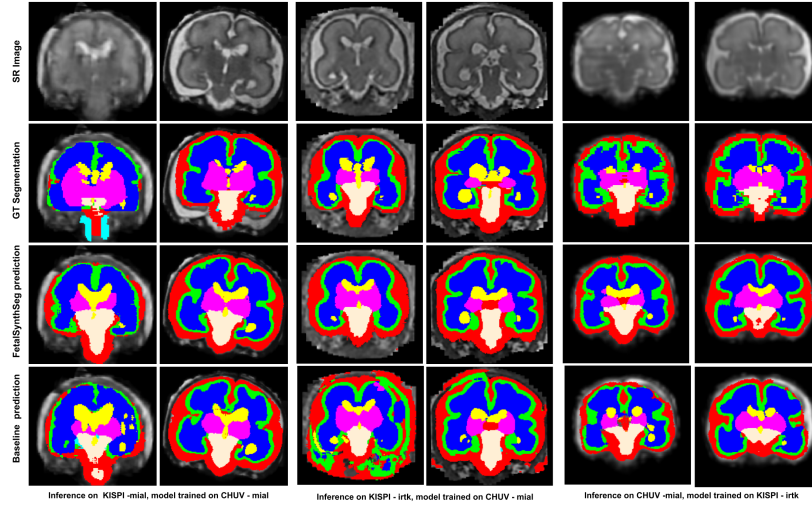| Synthetic Data Generator Hyperparameters | | | | | | |
|---|---|---|---|---|---|---|
| $a_{sc}$ | 0.9 | $a_{tr}$ | -10 | $r_{HR}$ | 0.5 | |
| $b_{sc}$ | 1.1 | $b_{tr}$ | 10 | $b_{res}$ | 0.5 | |
| Augmentations Hyperparameters$_{-Probability}$ | | | | | | |
| $\gamma_{range-0.5}$ | $0.5 - 1.5$ | $scale_{range-0.5}$ | $-0.1 - 0.1$ | $\sigma_{noise-0.5}$ | 0.1 | $\mu_{noise-0.5}$ 0 |
| $rotation_{range-0.5}$ | $-0.2 - 0.2$ | $shear_{range-0.5}$ | $-0.1 - 0.1$ | $\sigma_{smooth-0.7}$ | $0.5 - 1.5$ | |

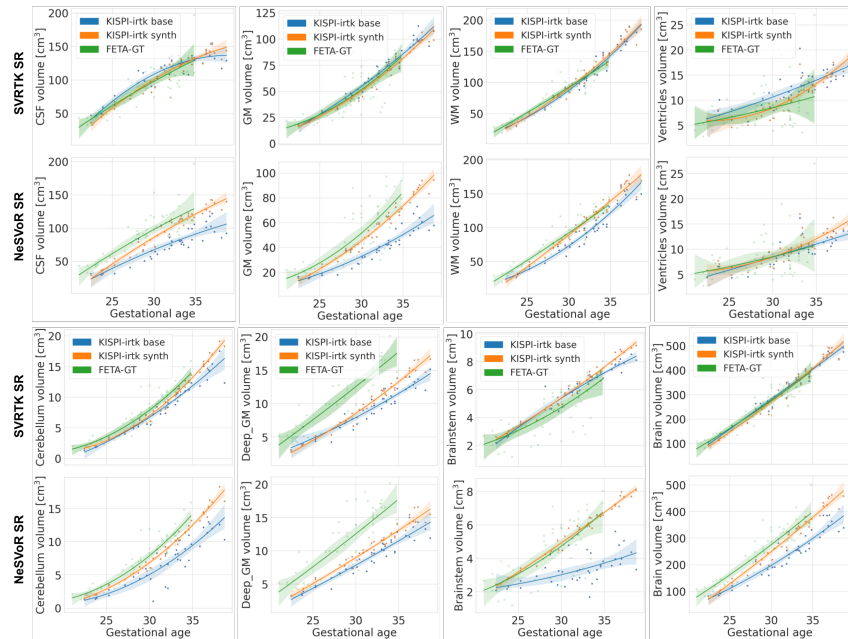**Fig. S2.** Cross-domain model inference qualitative results



**Fig. S3.** Segmented tissue volumes vs GA for KCL data for all tissues.