

NANOGP1, a tandem duplicate of *NANOG*, exhibits partial functional conservation in human naïve pluripotent stem cells

Katsiaryna Maskalenka^{1,*}, Gökberk Alagöz^{2,*}, Felix Krueger³, Joshua Wright¹, Maria Rostovskaya¹, Asif Nakhuda⁴, Adam Bendall¹, Christel Krueger¹, Simon Walker⁵, Aylwyn Scally² and Peter J. Rugg-Gunn^{1,6,7,‡}

ABSTRACT

Gene duplication events can drive evolution by providing genetic material for new gene functions, and they create opportunities for diverse developmental strategies to emerge between species. To study the contribution of duplicated genes to human early development, we examined the evolution and function of *NANOGP1*, a tandem duplicate of the transcription factor *NANOG*. We found that *NANOGP1* and *NANOG* have overlapping but distinct expression profiles, with high *NANOGP1* expression restricted to early epiblast cells and naïve-state pluripotent stem cells. Sequence analysis and epitope-tagging revealed that *NANOGP1* is protein coding with an intact homeobox domain. The duplication that created *NANOGP1* occurred earlier in primate evolution than previously thought and has been retained only in great apes, whereas Old World monkeys have disabled the gene in different ways, including homeodomain point mutations. *NANOGP1* is a strong inducer of naïve pluripotency; however, unlike *NANOG*, it is not required to maintain the undifferentiated status of human naïve pluripotent cells. By retaining expression, sequence and partial functional conservation with its ancestral copy, *NANOGP1* exemplifies how gene duplication and subfunctionalisation can contribute to transcription factor activity in human pluripotency and development.

KEY WORDS: Pluripotency, Reprogramming, Transcription factor, Gene duplication, Pseudogene, Evolution

INTRODUCTION

Gene duplication is an important driver of genome and species evolution. The majority of protein-coding genes and many non-coding regulatory sequences have arisen by duplication events (Magadum et al., 2013; Ohta, 2000). Most duplicated genes undergo functional decay due to silencing, loss-of-function mutations or lack of required regulatory regions (Magadum et al.,

2013). However, some duplicated genes are expressed, with the new copy either acquiring a novel function (neofunctionalisation) or sharing the ancestral function with the parental gene (subfunctionalisation). As a result, the emergence of a new copy of a gene or a regulatory sequence enables organisms to exploit new competitive advantages and to adapt to changing environments (Fares, 2014; Force et al., 1999; Kondrashov and Kondrashov, 2006).

Human evolution and development have been driven in many cases by the gain of low-copy repeats called segmental duplications. Over 5% of the human genome consists of segmental duplications, typically with more than 90% identity shared between the ancestral and the duplicated copies (Bailey et al., 2002; Marques-Bonet et al., 2009a). This percentage of duplicated regions is remarkably high compared to Old World monkeys, such as macaques, where only 1.5% of the genome consists of such duplicates (Marques-Bonet et al., 2009a). A burst of duplication events followed the divergence of apes from Old World monkeys, and these copies account for ~80% of modern, human-specific duplications (Marques-Bonet et al., 2009b). For example, two gene duplicates – *SRGAP2C* and *ARHGAP11B* – that are expressed in the developing human brain are proposed to have had a key role in the evolutionary expansion of the human neocortex (Charrier et al., 2012; Dennis and Eichler, 2016; Florio et al., 2015). However, the consequences of duplications underpinning such contributions remain largely undefined. Therefore, gene duplication events could be a major, unexplored driver of the divergence between mammalian developmental programmes, yet, for most duplicated genes, their contribution to these early developmental programmes is poorly understood.

The core pluripotency transcription factor *NANOG* has a high number of duplicated copies in the human genome, and could therefore serve as a paradigm for studying the impact of gene duplication events on early development. High expression levels of *NANOG* are crucial for maintaining the undifferentiated status of human naïve and primed states of pluripotency (Guo et al., 2021; Hyslop et al., 2005; Lie et al., 2012; Vallier et al., 2009; Zaehres et al., 2005). If any of its duplicated copies are also highly expressed, this would raise the possibility that they might have an unanticipated role in human pluripotent cells. Ten of the 11 duplicates of *NANOG* are processed pseudogenes (copies of mRNAs that have been reverse transcribed and inserted into the genome), which lack regulatory sequences and possess various mutations that have led to their functional decay (Booth and Holland, 2004). Only one member of the *NANOG* pseudogene family – *NANOGP1* – is unprocessed (Booth and Holland, 2004). *NANOGP1* transcripts are detected in leukaemia cells, adult testes and conventional or primed-state human pluripotent stem cells (hPSCs; naïve-state hPSCs have not been examined) (Eberle et al.,

¹Epigenetics Programme, Babraham Institute, Cambridge CB22 3AT, UK.

²Department of Genetics, University of Cambridge, Cambridge CB2 3EH, UK.

³Bioinformatics Group, Babraham Institute, Cambridge CB22 3AT, UK. ⁴Gene Targeting Facility, Babraham Institute, Cambridge CB22 3AT, UK. ⁵Imaging Facility, Babraham Institute, Cambridge CB22 3AT, UK. ⁶Wellcome-MRC Cambridge Stem Cell Institute, Cambridge CB2 0AW, UK. ⁷Centre for Trophoblast Research, University of Cambridge, Cambridge CB2 3EG, UK.

*These authors contributed equally to this work

‡Author for correspondence (peter.rugg-gunn@babraham.ac.uk)

 P.J.R.-G., 0000-0002-9601-5949

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0>), which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

Handling Editor: Maria Elena Torres-Padilla
Received 26 July 2022; Accepted 16 December 2022

2010; Hart et al., 2004). *NANOG* and *NANOGP1* share 97% coding region homology and have a similar exon-intron structure, suggesting that *NANOGP1* has probably undergone selection-driven conservation (Booth and Holland, 2004; Fairbanks and Maughan, 2006). Previous studies have reached contradictory conclusions about whether *NANOGP1* encodes a full-length protein (Booth and Holland, 2004; Eberle et al., 2010). If *NANOGP1* uses the equivalent translation initiation codon as *NANOG*, then, owing to a base pair substitution, the resultant protein would contain only the first eight amino acid residues. However, *NANOGP1* could use an alternative downstream initiation start codon that would encode a near full-length protein. This predicted *NANOGP1* protein, if expressed, would have an intact homeodomain and transactivation domain, which are responsible for the protein dimerisation, DNA binding and pluripotency maintenance functions of *NANOG* and its orthologs (Chambers et al., 2003; Chang et al., 2009; Hart et al., 2004; Mullin et al., 2021; Oh et al., 2005; Theunissen et al., 2011). Whether endogenous *NANOGP1* is translated into a protein has not been determined. This uncertainty about the predicted *NANOGP1* open reading frame led to the belief that *NANOGP1* does not encode a protein (Booth and Holland, 2004), and *NANOGP1* is currently classified as a non-protein-encoding pseudogene in the Ensembl repository.

Because *NANOG* has a central role in regulating pluripotency, it is important to establish whether *NANOGP1* is a protein-coding gene that could also have functional capabilities. Here, we show that the *NANOGP1* protein is expressed in naïve-state hPSCs. We determined that *NANOG* and *NANOGP1* have overlapping but not identical expression patterns in human embryos and stem cell lines. We found that, in contrast to *NANOG*, *NANOGP1* is not required to maintain undifferentiated naïve hPSCs, but *NANOGP1* can fulfil other functional roles of *NANOG*, including reprogramming and autorepressive activities. Furthermore, genetic analysis established that the duplication by which *NANOGP1* was formed occurred earlier than previously thought and before the divergence of apes and Old World monkeys, and that the gene has been decayed in Old World monkeys but retained in great apes. By establishing that *NANOGP1* has retained partial functional conservation with its ancestral copy *NANOG*, our study sheds light on the role of gene duplication and subfunctionalisation on human pluripotency and development.

RESULTS

Identification of pseudogenes, including *NANOGP1*, that are highly expressed in human naïve pluripotent stem cells

To investigate pseudogene expression in human pluripotent cells, we first analysed transcript levels of pseudogenes in naïve-state hPSCs using RNA sequencing. We selected 1880 protein-coding genes in the human genome that have pseudogene copies (totalling 6922 transcripts; Ensembl 104 annotation). Overall, 486 pseudogenes were detected with an expression value of $\log_2\text{RPM} > 0$ in naïve hPSCs (Fig. 1A). Highly expressed pseudogenes have ancestral genes that are enriched for roles in RNA binding and translation (Fig. S1A), and have higher sequence conservation when compared with pseudogenes expressed at lower levels (Fig. S1B). We also found that several key pluripotency factors, including *NANOG*, *POU5F1* (also known as *OCT4*) and *DPPA3*, had highly expressed pseudogenes in naïve hPSCs (Fig. 1B, Fig. S1C-E). Four of these duplicated genes – *NANOGP1*, *POU5F1P4*, *POU5F1P3* and *DPPA3P2* – were within the top ~2% of all pseudogenes ranked by expression levels and their levels approached those of their ancestral copies

(Fig. 1B, Fig. S1C-E). In addition to the duplicated pseudogene *NANOGP1* that was highly expressed, the processed and truncated gene *NANOGP8* also had a substantial number of mapped reads (Fig. S1C). *POU5F1P4*, *POU5F1P3*, *DPPA3P2* and *NANOGP8* are processed copies, whereas *NANOGP1* was of specific interest because it has been formed by tandem duplication, is unprocessed and is located in the same locus as its ancestral copy, *NANOG*. Together, these results uncover the large set of pseudogenes that are expressed in naïve hPSCs. In particular, the high expression of the duplicated pseudogene *NANOGP1* raises the possibility that this gene might have an unanticipated role in human pluripotent cells.

NANOG and *NANOGP1* have overlapping but distinct expression patterns

To study the expression pattern of *NANOGP1*, we next compared RNA-seq datasets of naïve and primed hPSCs (Collier et al., 2017), which are cell types that correspond to early and late epiblast cells of the human embryo, respectively. Although *NANOGP1* is a duplicated copy of *NANOG*, there were sufficient sequence differences between the transcripts of the two genes to uniquely assign RNA-seq reads to each gene (Sequence Divergence Rate of 0.013) (Fig. S2). We also confirmed that *NANOG* reads do not map to the *NANOGP1* locus, and vice versa, when using a high mapping quality value ($\text{MAPQ} > 20$). The transcriptional analysis revealed notable differences in the expression patterns of *NANOG* and *NANOGP1*. Whereas *NANOG* is highly expressed in both naïve and primed hPSCs, *NANOGP1* is highly expressed in only naïve hPSCs, and is substantially downregulated in primed hPSCs (Fig. 1C). (Previous studies examined only primed hPSCs.) This finding was extended by analysing multiple RNA-seq data sets of different naïve and primed hPSC lines, including embryo-derived and reprogrammed cell lines, and cultured in different media conditions (Fig. 1D) (Guo et al., 2016; Pastor et al., 2016; Takashima et al., 2014; Theunissen et al., 2016).

To test whether the distinct expression patterns are also observed *in vivo*, we reanalysed single-cell RNA-seq (scRNA-seq) datasets from human embryos (Petropoulos et al., 2016; Xiang et al., 2020). Like *NANOG*, *NANOGP1* was highly expressed in epiblast but not trophectoderm lineages (Fig. 1E). *NANOG* and *NANOGP1* expression was well-correlated in pre-implantation epiblast cells (Fig. S3A). Interestingly, we found that *NANOGP1* might be expressed in a subpopulation of primitive endoderm cells, although available cell numbers are low for this lineage (Fig. 1E). *NANOGP1* and *NANOG* transcripts were abundant throughout epiblast development, up until day 14, at which point *NANOGP1* levels were abruptly reduced (Fig. 1F). In contrast, *NANOG* expression levels remained high including on day 14 (Fig. 1F). This developmental expression pattern therefore mirrored the state-specific differences between naïve and primed hPSCs, further confirming the overlapping but distinct expression profiles of the two genes. Finally, as *NANOG* is expressed in germ cells, we examined published RNA-seq data of *in vivo* germ cells (Gkoutela et al., 2015) and found that *NANOGP1* transcripts are also detected at high levels that are comparable with *NANOG* (Fig. S3B). Overall, these results show that *NANOGP1* is dynamically expressed in hPSCs and developing human embryos, which is an expression pattern that suggests a conserved potential role for *NANOGP1* in human early development.

NANOGP1 transcript and protein isoform sequences are highly similar to those of *NANOG*

The high expression and sequence read coverage of *NANOGP1* in naïve hPSCs enabled us to examine its mRNA structure, splicing

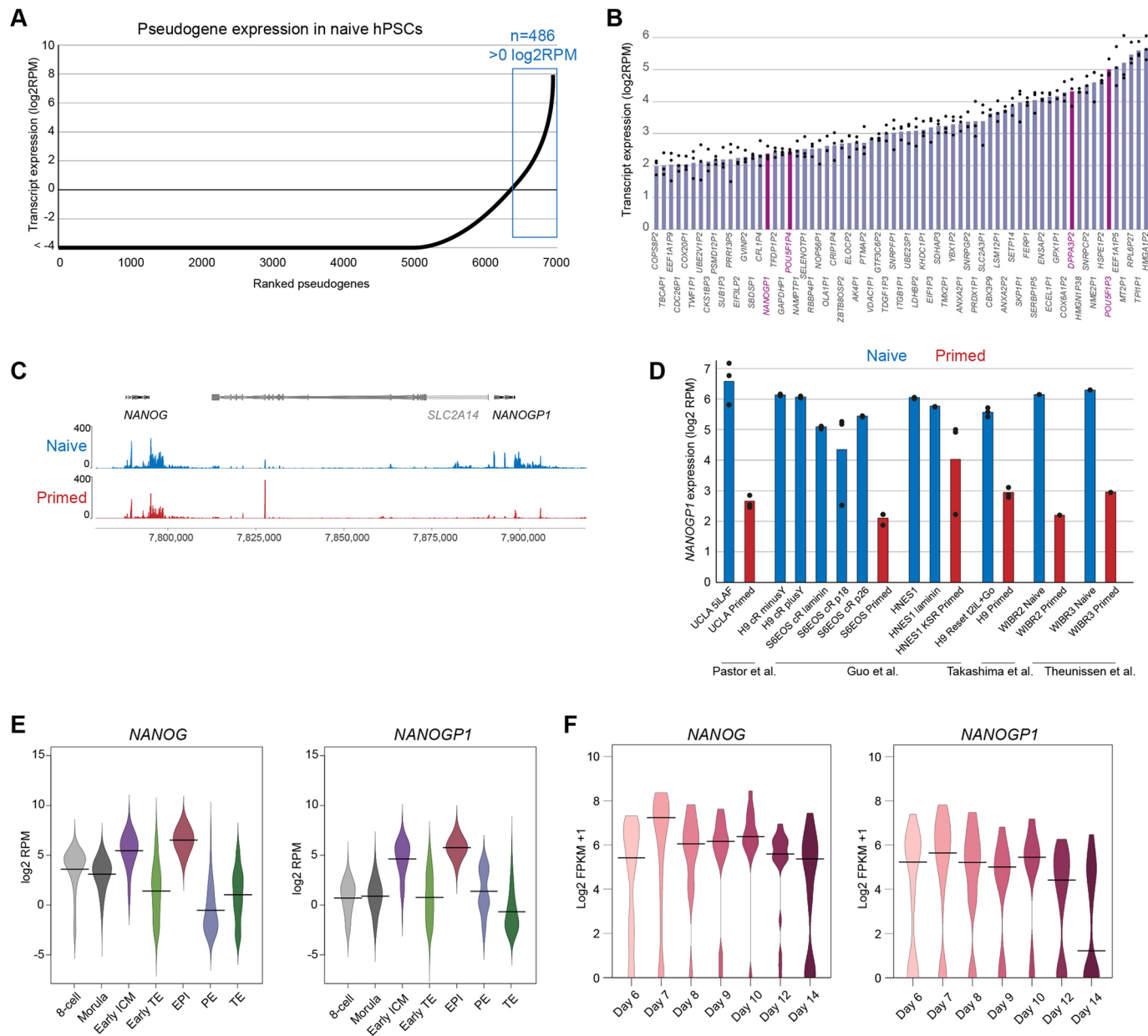


Fig. 1. *NANOGP1* is highly expressed in human naïve pluripotent stem cells and epiblast cells. (A) Ranked expression of 6922 pseudogenes in naïve hPSCs. Analysis was performed using a custom annotation of pseudogenes. The y-axis has been cut off at $-4 \log_2$ RPM. (B) Examples of highly expressed pseudogenes in naïve hPSCs. Pseudogenes of pluripotency factors are in dark purple. Analysis performed using a custom annotation of pseudogenes. Data show mean and data points from three independent samples. (C) RNA-seq data for *NANOG*, *SLC2A14* and *NANOGP1* in naïve and primed hPSCs (Collier et al., 2017). (D) *NANOGP1* expression in naïve (blue) and primed (red) hPSC lines (Guo et al., 2016; Pastor et al., 2016; Takashima et al., 2014; Theunissen et al., 2016). Data show mean and data points from three independent samples (except for the WIBR2 and WIBR3 lines, which have one data point each). (E) *NANOG* and *NANOGP1* expression in human pre-implantation embryos (Petropoulos et al., 2016). 8 cell, eight-cell stage ($n=78$); Morula ($n=185$); early ICM, early inner cell mass ($n=66$); early TE, early trophectoderm ($n=227$); EPI, epiblast ($n=45$); PE, primitive endoderm ($n=30$); TE, trophectoderm ($n=715$). Horizontal lines indicate the median. (F) *NANOG* and *NANOGP1* expression in epiblast cells from human peri-implantation and early post-implantation cultured embryos (Xiang et al., 2020). Day 6 ($n=60$); day 7 ($n=33$); day 8 ($n=11$); day 9 ($n=12$); day 10 ($n=14$); day 12 ($n=22$); day 14 ($n=26$). Horizontal lines indicate the median.

patterns and open reading frame sequences. Using published RNA-seq data (Takashima et al., 2014), this analysis identified three *NANOGP1* mRNA isoforms that differed due to alternative splicing between exons 3 and 4 (Fig. 2A). This pattern was consistent in additional naïve hPSC lines (Fig. S4) (Theunissen et al., 2016; Pastor et al., 2016). No splicing to a putative upstream exon was detected, as had been previously considered (Booth and Holland, 2004). According to the splicing analysis in our study, the first

NANOGP1 exon was the same as that of *NANOG*. Owing to a point mutation within exon 1, the most likely translation initiation codon for *NANOGP1* is 117 bp downstream of the equivalent initiation codon used by *NANOG* (Fig. 2B). This results in the open reading frame of *NANOGP1* lacking the first 39 amino acids compared with *NANOG* (Fig. 2C), which is a finding that is consistent with earlier predictions (Booth and Holland, 2004; Hart et al., 2004). Outside the first exon, the sequences encoding the main protein domains

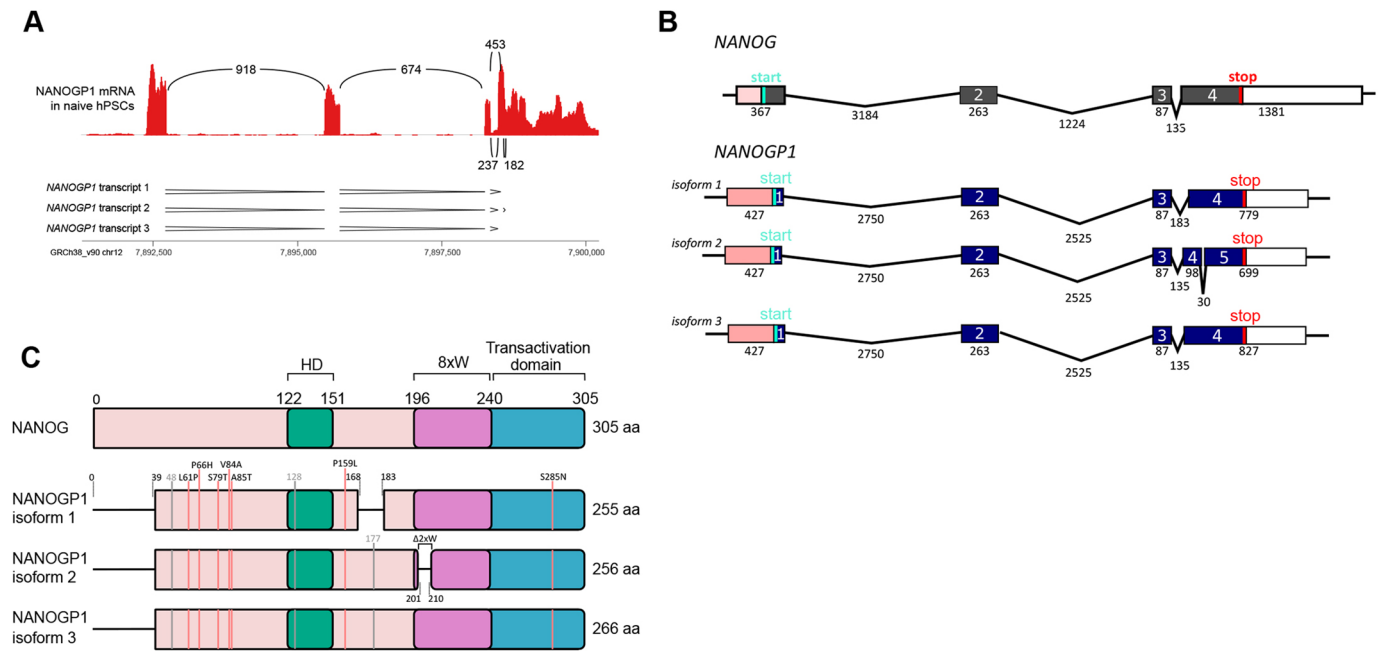


Fig. 2. Predicted open reading frame structure of *NANOGP1*. (A) Splicing analysis of *NANOGP1* in naive hPSCs (Takashima et al., 2014). The numbers in between the RNA-seq peaks indicate the number of times a splicing was measured. The three predicted patterns of transcript splicing are underneath. (B) Predicted transcript isoforms of *NANOGP1*, including the size of exons and introns (in bp), and translation start and stop codons. The transcript structure of *NANOG* is shown for comparison. (C) Predicted *NANOGP1* open reading frame (ORF) variants and domain structures. The ORF of *NANOG* is shown for comparison. Differences in the *NANOGP1* ORFs versus the *NANOG* ORF are indicated. Amino acid substitutions caused by missense DNA changes are labelled by red vertical lines; silent changes are labelled by grey vertical lines. 8xW, tryptophan-rich subdomain/region containing eight tryptophan (W) residues; $\Delta 2 \times W$, deletion of two tryptophan residues from the tryptophan-rich subdomain; HD, DNA-binding homeodomain.

of *NANOG*, including the homeobox domain, tryptophan repeats and C-terminal transactivation domain, were all present and fully conserved in all the predicted *NANOGP1* open reading frames and isoforms (Fig. 2C). Several point mutations and two smaller deletions in isoforms 1 and 2 were detected outside the main domains (Fig. 2C). Overall, these results show that the predicted sequences, exon structures and functional domains of *NANOGP1* are very similar to *NANOG*.

***NANOGP1* gene and protein sequences are highly conserved in great apes**

We next examined the boundaries of the *NANOG/NANOGP1* duplication in the human genome. We self-aligned a 250 kb region containing *NANOG*, *NANOGP1*, *SLC2A14*, *SLC2A3* and *NANOGNB*, plus their flanking regions on both sides (Fig. 3A). Three large domains of duplication were identified: (1) *NANOG* and *NANOGP1*; (2) *SLC2A14* and *SLC2A3*; and (3) a *SLC2A3* downstream region (Fig. 3A,B). These results are consistent with a duplication event that involved copying and inserting an ~80 kb region containing *NANOG* and *SLC2A14* into a new location immediately downstream of its original position, and which resulted in the formation of the *NANOG/NANOGP1* duplication.

To better understand the origins and conservation of the *NANOG/NANOGP1* duplication, we manually examined gene lengths, genomic positions and gene orientation data from genome assemblies of non-human apes, Old and New World monkeys, and prosimians. We searched for unambiguous matches to *NANOGP1* in each assembly and annotated it where present, as this annotation was absent from most of the non-human genomes. We then aligned identified *NANOGP1* sequences to their corresponding *NANOG* counterparts (Fig. S5A,B). Our analysis revealed that the *NANOGP1* sequence is present in some ape and

Old World monkey genomes, but not in New World monkey or prosimian genomes (Fig. 3C, Fig. S5A). This finding suggests that the duplication event occurred before the split between apes and Old World monkeys (30-35 million years ago, Mya) but more recently than the split between the Old World and New World monkeys (40-50 Mya) (Pozzi et al., 2014), and was followed by full or partial deletion on some lineages outside the great apes (Fig. S5A-C). We note, however, that the marmoset genome (New World monkey) contains *SLC2A3*, which is a duplicated gene of *SLC2A14* (Fig. 3C). An alternative interpretation, therefore, is that the duplication event pre-dated ~50 Mya and that *NANOGP1* was subsequently lost from the marmoset genome, or that there were two separate duplication events: the first for *SLC2A14/SLC2A3* and the second for *NANOG/NANOGP1*.

NANOGP1 sequences are present in most of the examined Old World monkey and ape species (Fig. 3C). Interestingly, however, an intact copy of *NANOGP1* is present only in great apes and, instead, the other species have inactivated *NANOGP1* in different ways. Some species, such as gibbon, have deleted the entire gene, whereas others, including the green monkey and crab-eating macaque, have partial deletions of *NANOGP1* (Fig. 3C, Fig. S5A-C). These species have retained *SLC2A3*. Other species appeared initially to have retained intact *NANOGP1*, but closer inspection uncovered small, critical mutations that are predicted to disable the protein. For example, *Rhesus macaque* contains a full-length *NANOGP1* sequence, but crucially has a non-synonymous amino acid change within the homeodomain (Fig. 3D). The affected amino acid, M54I, confers the DNA-binding specificity of *NANOG* (Weiler et al., 1998). The likely consequence of this change is altered target sequence recognition because the homeobox protein PBX1, which also has an isoleucine at position 54, has a consensus motif of TGAT that differs from the canonical TAAT motif of *NANOG* (Chang

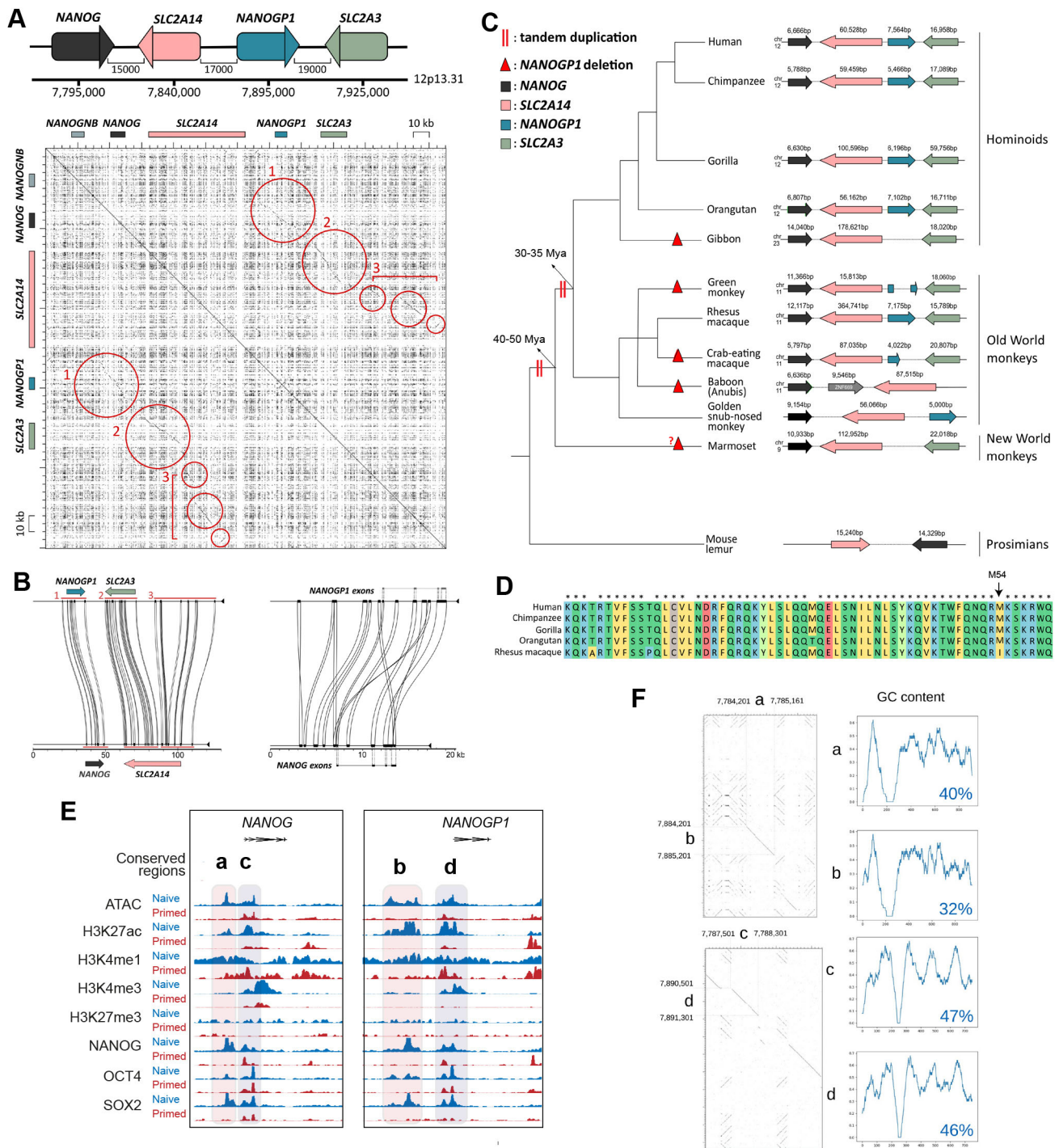


Fig. 3. *NANOGP1* duplication in human evolution. (A) Top: *NANOG/NANOGP1* tandem duplication locus [distance (bp) between the genes/pseudogene]. Bottom: self-alignment of a 250 kb region containing *NANOGNB*, *NANOG*, *NANOGP1* and another duplicated gene pair, *SLC2A14* and *SLC2A3* (genes indicated by boxes along x- and y-axes). Individual dots represent matching base pairs between the two aligned sequences. Circles indicate three areas of high sequence conservation between the ancestral and duplicated regions. (B) Sequence similarity and locations of the three regions identified in A (left) and between the exons and upstream regions of *NANOG* and *NANOGP1* (right). (C) Conservation of the *NANOG/NANOGP1* tandem duplication locus across species. Predicted duplication dates are indicated with two red vertical lines; predicted *NANOGP1* deletion events are indicated with red triangles. (D) Amino acid alignment compares the homeodomain sequences of *NANOGP1* orthologs. Colour indicates different types of amino acids, according to their biochemical properties. Asterisks indicate that the amino acid is the same for all aligned sequences. (E) ATAC-seq (Pastor et al., 2016) and ChIP-seq (Chovanec et al., 2021) profiles across the *NANOG* and *NANOGP1* loci in naïve and primed hPSCs. The sequences labelled 'a-d' indicate two duplicated pairs of regulatory regions. (F) Comparison of the regulatory regions a-d. Left: individual dots represent matching base pairs between the two aligned sequences. Right: GC content ratio graphs in which the x-axis represents the length of a putative regulatory region in bp, and the y-axis shows GC content within 30 bp sliding windows. The average GC content ratios over the indicated regions are shown.

et al., 1996; Piper et al., 1999). The function of *NANOGP1* in *Rhesus macaque* is therefore likely to be compromised. In contrast, the homeodomain sequences are intact for *NANOGP1* in human, chimpanzee and gorilla (Fig. 3D).

These results show that a duplication event around 40 Mya created the *NANOG/NANOGP1* duplicated region that is present in the genomes of Old World monkeys and apes. *NANOGP1* has subsequently been disabled in most of the primate genomes via different alterations. Great apes, however, have retained intact gene and protein sequences, suggesting the potential presence of evolutionary pressure to maintain *NANOGP1* in those species.

Putative regulatory regions upstream of *NANOGP1* were formed in the tandem duplication event

In addition to highly conserved exons, we also found distal regions that were conserved. Examining the sequence conservation and chromatin marks at the *NANOG/NANOGP1* locus revealed the location of several putative regulatory regions that overlapped with elements previously annotated as enhancers and super-enhancers (Fig. 3E, Fig. S6) (Chovanec et al., 2021). Six of these regions were identified near to *NANOGP1*, and four were positioned as two pairs directly upstream of *NANOG* (a, c) and *NANOGP1* (b, d) (Fig. 3E, Fig. S6). Pairwise alignments showed that the sequences within the two individual pairs, a/b and c/d, were similar; additionally, each pair had matching GC content profiles, providing further evidence that they had formed from a duplication event (Fig. 3F). For the c/d pair, the GC content ratios were close to typical GC content ratio values that average ~50% in promoter regions (Villar et al., 2015), in contrast to the a/b pair that had lower GC content values (Fig. 3F). Together with the chromatin profiles, such as the promoter-associated modification H3K4me3, this allowed us to conclude that c/d are likely serve as promoters and a/b as enhancers.

According to ATAC-seq profiles (Pastor et al., 2018), sites a, b, c and d have highly accessible chromatin (Fig. 3E). Additionally, all four regions have high levels of active histone modifications – H3K27ac, H3K4me1 and H3K4me3 – and are bound by pluripotency factors in either one or both hPSC states (Fig. 3E) (Chovanec et al., 2021). The putative promoters c and d appeared active in both naïve and primed hPSC states, and were hence referred to as ‘shared’, while the putative enhancers a and b were predominantly marked as active in the naïve hPSCs. The pattern of transcription factor occupancy and chromatin annotations were similar for *NANOG* and *NANOGP1* at their putative promoter regions. The only prominent differences were for SOX2 and H3K4me3 levels within the shared putative promoters, where SOX2 and H3K4me3 peaks were detected near to *NANOG* in both primed and naïve hPSCs, but were present only in naïve hPSCs at the *NANOGP1* locus.

These results demonstrate that *NANOGP1* is integrated within the regulatory circuitry of pluripotent cells through OCT4, SOX2 and NANOG binding. The similarities in enhancer conservation and annotations could also help to explain the overlap of *NANOGP1* and *NANOG* expression patterns in human embryos and naïve hPSCs, and differences at the *NANOGP1* promoter in primed hPSCs correlate with reduced *NANOGP1* expression in those cells.

NANOGP1 encodes a protein that is expressed in naïve pluripotent stem cells

Although *NANOGP1* is currently annotated as a non-protein-encoding pseudogene, our revised sequence analysis suggested that the transcript should encode a protein of at least 255 amino acids. We therefore sought to establish whether NANOGP1 protein is

detectable in naïve hPSCs. The close similarity in the predicted protein sequences of NANOGP1 and NANOG means there are no antibodies to detect NANOGP1 only, so we chose to insert an epitope tag into the endogenous *NANOGP1*-coding sequence through homology directed repair (HDR). Pilot experiments established that the most efficient *in vivo* DNA cutting efficiency was obtained with CRISPR-Cas12a endonuclease targeting near to the start codon of *NANOGP1* (Table S1).

We therefore used Cas12a ribonucleoprotein (RNP) and single stranded DNA (ssDNA) templates to insert V5 and 3xFLAG epitope tags into the endogenous *NANOGP1*-coding sequence in naïve hPSCs (Fig. 4A,B). We detected nuclear-localised expression of epitope-tagged NANOGP1 in polyclonal naïve hPSCs by immunostaining (Fig. 4C). Epitope-tagged NANOGP1 was also identified after immunoprecipitation and western blotting (Fig. 4D). The specificity of the epitope-tagged protein was confirmed by using two different anti-NANOG antibodies for the western blot: one that recognises the C termini of NANOG and NANOGP1, and one that recognises the N terminus of NANOG but not NANOGP1 (owing to the N-terminal truncation of NANOGP1). These results establish that, in contrast to current annotations, *NANOGP1* is a protein-coding gene and its product is expressed in naïve hPSCs.

The discovery of NANOGP1 protein in naïve hPSCs prompted us to investigate whether this factor might have functional roles in naïve pluripotency. *NANOG* has several known functions in naïve pluripotent stem cells, including (1) a gene autorepressive ability that was identified in mouse pluripotent stem cells (Navarro et al., 2012), (2) suppression of the transcription of the trophoblast marker genes *GATA2*, *GATA3* and *TFAP2C* (Guo et al., 2021), and (3) an ability to reprogramme primed hPSCs towards the naïve state when overexpressed together with *KLF2* (Takashima et al., 2014; Theunissen et al., 2014). These three aspects of *NANOG* function were tested in relation to *NANOGP1* in the following sections.

NANOGP1 has repressive activity on *NANOG* and *NANOGP1*

Ectopic *Nanog* overexpression in mouse pluripotent stem cells leads to the autorepression of endogenous *Nanog* expression (Navarro et al., 2012). To test whether *NANOG* and/or *NANOGP1* overexpression has a similar effect in human naïve pluripotency, we established hPSC lines containing doxycycline-inducible *NANOG* and *NANOGP1* transgenes (Fig. 5A,B). The induction of *NANOG* expression led to the downregulation of endogenous *NANOG* (Fig. 5C), thereby establishing that, as for mouse, human *NANOG* also has gene autorepressive activity. Interestingly, endogenous *NANOGP1* was also downregulated (Fig. 5C). Importantly, the overexpression of *NANOGP1* also suppressed the expression of *NANOG* and endogenous *NANOGP1* (Fig. 5D). This effect was also observed in primed hPSCs (Fig. 5E). These results establish that *NANOGP1* has a conserved autorepressive function.

NANOGP1 can reprogramme human primed pluripotent stem cells into a naïve state

The short-term, enforced expression of *NANOG* and *KLF2* facilitates the reprogramming of primed hPSCs into the naïve state (Takashima et al., 2014; Theunissen et al., 2014). We therefore investigated whether *NANOGP1* is also capable of promoting primed to naïve reprogramming, to ascertain whether *NANOGP1* can fulfil the role of *NANOG* in a direct functional test. *NANOGP1* was overexpressed together with *KLF2* in primed hPSCs using a doxycycline-inducible system in minimal 2i+LIF medium (Fig. 6A). We tested all three *NANOGP1* isoforms separately. To monitor and select for transgene expression, *NANOGP1* was

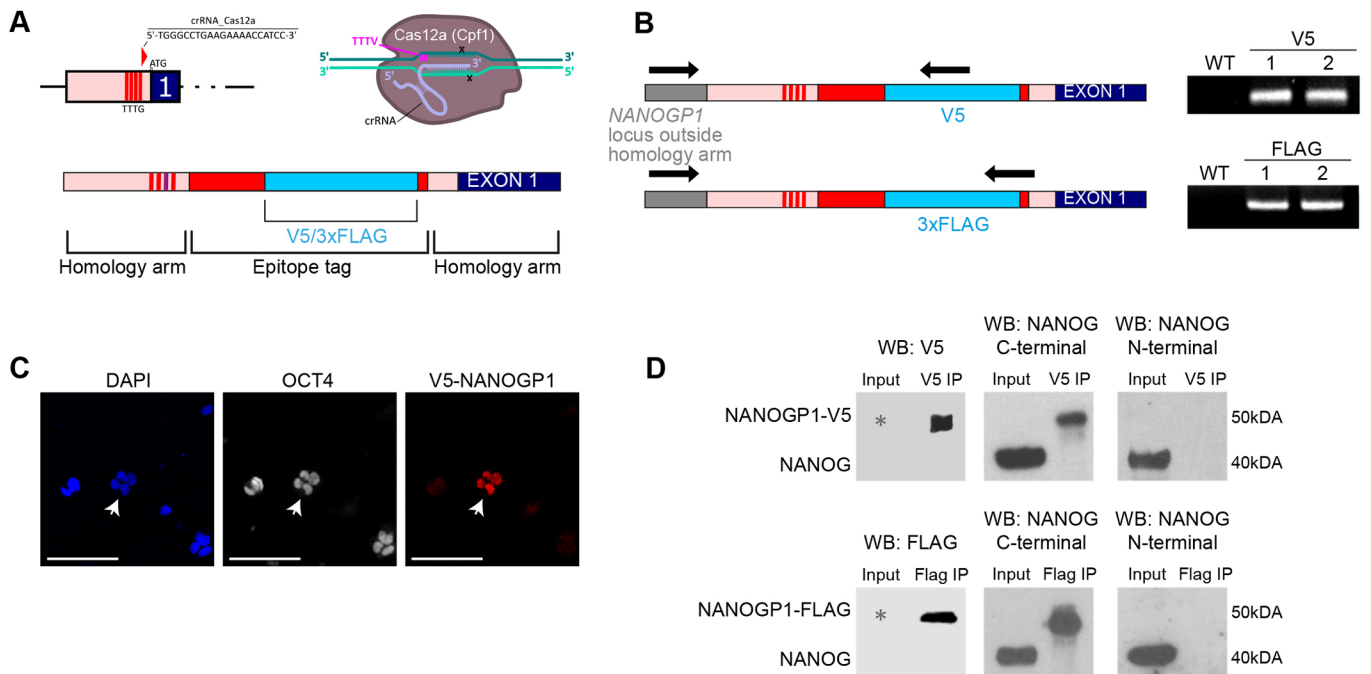


Fig. 4. *NANOGP1* encodes a protein that is expressed in human pluripotent cells. (A) CRISPR/Cas12a strategy to target *NANOGP1* and insert in-frame V5 or 3xFLAG epitope tags. (B) Left: genotyping strategy with primer positions (arrows). Right: integration of the tags into the *NANOGP1* locus in naïve hPSCs. WT, untransfected naïve hPSCs; V5 lane 1 and V5 lane 2, two independent lines with V5 integrated at the *NANOGP1* locus; FLAG lane 1 and FLAG lane 2, two independent lines with 3xFLAG integrated at the *NANOGP1* locus. (C) Nuclear localisation of V5-NANOGP1 in small colonies of polyclonal transgenic naïve hPSCs, and overlap with OCT4 and DAPI signal. White arrows indicate the V5-positive colony. The other visible colonies are V5 negative and presumably not successfully targeted. Scale bars: 100 μ m. (D) Western blot of co-immunoprecipitation experiments. Protein samples from transgenic polyclonal naïve hPSCs were immunoprecipitated with either V5 (upper) or FLAG (lower) antibodies. The immunoprecipitated material was examined by western blot using antibodies against the epitope tag (left), the NANOG C terminus that also detects NANOGP1 (centre), and the NANOG N terminus that does not detect NANOGP1 due to an N-terminal deletion (right). The grey asterisks indicate that, due to the low number of NANOGP1-epitope tagged cells in the polyclonal population, the proteins were detected only in the immunoprecipitated samples and not in the input samples.

co-translated with *GFP* via an internal ribosome entry site, and *KLF2* with *RFP*. Before reprogramming, we ensured comparable overexpression levels in all lines by inducing the cells with doxycycline for 24 h and flow sorting the appropriate *GFP*^{+*RFP*⁺ or *RFP*-only⁺ cell populations (Fig. S7A). The following day, the cells were switched to 2i+LIF medium with doxycycline to initiate reprogramming.}

By day 12 of reprogramming in these conditions, we observed numerous domed colonies with naïve hPSC morphology in the *NANOGP1*+*KLF2* cultures. The cells had upregulated naïve pluripotency markers, including *DPPA3* and *TFCP2L1*, and maintained high *POU5F1* expression (Fig. 6B). All three *NANOGP1* isoforms showed similar effects. These changes were comparable with the positive control cells expressing *NANOG* and *KLF2*. The reprogrammed colonies were positive for alkaline phosphatase activity, and the number of positive colonies was similar when comparing cultures overexpressing either *NANOGP1* or *NANOG* (Fig. 6C, Fig. S7B). Flow cytometry analysis using cell-surface markers of naïve pluripotency (CD24 negative, CD75 positive and *SUSD2* positive) (Bredenkamp et al., 2019a; Collier et al., 2017; Shakiba et al., 2015; Wojdyla et al., 2020) validated successful pluripotent state conversion in the *NANOGP1*-overexpressing cells (Fig. 6D,E). Importantly, in all of the assays, the overexpression of *KLF2* alone did not induce reprogramming, confirming the crucial contribution of *NANOGP1* in establishing naïve pluripotency. The change in pluripotent state was stable because the *NANOGP1*-induced reprogrammed cells retained their cell-surface marker phenotype when cultured for seven passages

without doxycycline (Fig. 6F). Overall, these results lead us to conclude that, like *NANOG*, *NANOGP1* is capable of reprogramming hPSCs into the naïve state, thereby demonstrating functional conservation in igniting the naïve pluripotency network.

As *NANOGP1* levels are substantially lower in primed cells compared with naïve cells, we examined whether enforced *NANOGP1* expression can disrupt the transition of naïve hPSCs into a primed state. We used doxycycline to induce *NANOGP1* expression in naïve hPSCs and immediately changed the conditions to promote naïve-to-primed capacitation in the presence of doxycycline (Fig. 6G) (Rostovskaya et al., 2019). After 6 days, flow cytometry analysis revealed a strong reduction in CD24/SSEA4 double-positive primed cells in *NANOGP1*-expressing conditions (Fig. 6H). In addition, the expression of *DUSP6*, a primed marker, was significantly reduced in *NANOGP1*-expressing cells compared with non-induced control cells (Fig. 6I). *TFCP2L1*, a naïve marker, decreased moderately after 1 day, but then failed to further decrease over the next 5 days, resulting in elevated levels compared with control cells (Fig. 6I). These findings suggest that *NANOGP1* downregulation might be needed for effective naïve-to-primed transition. However, *NANOG* was rapidly downregulated after *NANOGP1* induction (Fig. 6I), consistent with the repressive effect of *NANOGP1* on *NANOG*, and this effect could partly explain the defect in transitioning to a primed state. We also observed substantial cell death in *NANOGP1*-expressing cells at the later stages of capacitation (Fig. 6I). Taken together, these results establish that enforced expression of *NANOGP1* disrupts naïve to primed capacitation.

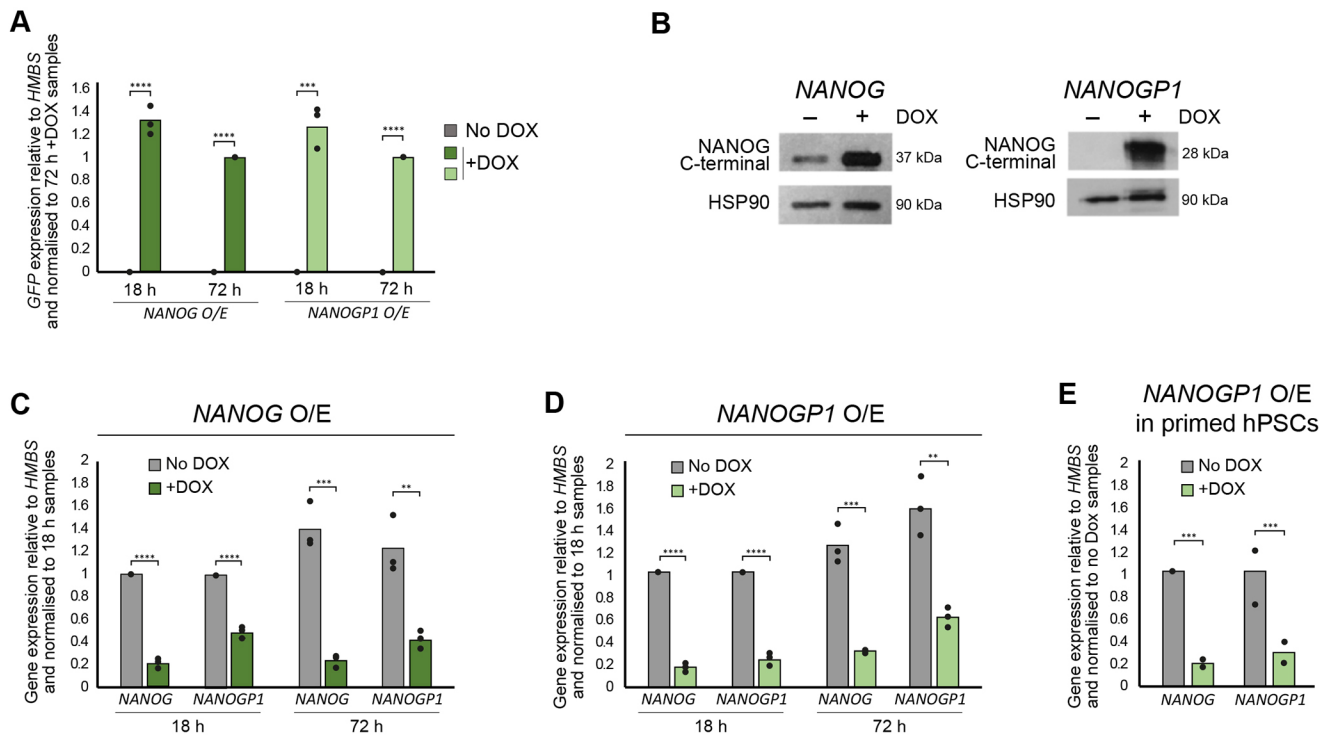


Fig. 5. NANOGP1 has gene autorepressive activity. (A) Induction of *NANOG*-GFP and *NANOGP1*-GFP transgenes in naïve hPSCs, as monitored by GFP expression. RT-qPCR values are relative to *HMBS* expression and normalised to the 72 h +DOX samples. Mean and data points from three independent samples are shown. Unpaired *t*-test (two-tailed; *** $P=0.0003$, **** $P<0.0001$). (B) Western blot showing DOX-induced overexpression of *NANOG* and *NANOGP1* in naïve hPSCs with DOX-inducible *NANOG* (C) and *NANOGP1* (D) transgenes. Primers target the 5' UTR of either *NANOG* or *NANOGP1*. RT-qPCR values are relative to *HMBS* expression and normalised to the 18 h samples. Mean and data points from three independent samples are shown. Unpaired *t*-test (two-tailed; ** $P<0.01$, *** $P<0.001$, **** $P<0.0001$). (E) Endogenous *NANOG* and *NANOGP1* expression levels in primed hPSCs with DOX-inducible *NANOGP1* transgene. Mean and data points from two independent samples are shown. Unpaired *t*-test (two-tailed; *** $P<0.001$). O/E, overexpression.

NANOGP1* is not required to maintain naïve pluripotency, unlike *NANOG

We next investigated whether *NANOGP1* supports the maintenance of human naïve pluripotency. A recent study showed that polyclonal cultures of *NANOG*-deficient naïve hPSCs upregulate several trophoblast lineage marker genes, thereby uncovering a potentially crucial role for *NANOG* in maintaining naïve pluripotency (Guo et al., 2021). However, the dynamics of the transcriptional response after *NANOG* perturbation, and the effect on gene expression programmes, have not been examined. We first aimed at better defining this important phenotype, which would also provide a suitable comparison for studying whether the loss of *NANOGP1* might show similar effects.

We established naïve hPSC lines expressing doxycycline-inducible CRISPRi (dCas9-KRAB) (Mandegar et al., 2016) that targeted the promoters of either *NANOG* or *NANOGP1* by gene-specific gRNAs (Fig. 7A). Treating the transgenic naïve hPSC lines with doxycycline caused the efficient and gene-specific knockdown of *NANOG* transcripts by 80%, and *NANOGP1* levels by 90% (Fig. 7B). *NANOG* protein was also strongly reduced after doxycycline treatment (Fig. 7C).

CRISPRi-mediated *NANOG* downregulation caused the naïve cells to lose their characteristic domed morphology and to visibly differentiate (Fig. 7D). Consistent with this, RNA-seq profiling over a 9-day time course revealed a strong transcriptional downregulation of naïve and core pluripotency factors (Fig. 7E, Fig. S8). Transcriptionally upregulated genes were associated strongly with the trophoblast lineage, including *GATA2*, *GATA3*, *CDX2*,

ESRRB and *TACSTD2*, and their induction was detected on day 2 and continued to increase in their expression up to day 9 (Fig. 7E, Fig. S8). Other categories associated with upregulated genes over the timecourse included processes associated with mesoderm cell types, and Hippo and Wnt signalling pathways (Fig. S8).

In contrast, the downregulation of *NANOGP1* did not cause naïve hPSCs to induce the expression of trophoblast marker genes or to change their morphology (Fig. 7D,E). Expression of pluripotency genes were unaltered (Fig. 7E) and, overall, far fewer differentially expressed genes were detected after *NANOGP1* downregulation compared with *NANOG* downregulation (Fig. 7F). The transcriptional responses after the knockdown of *NANOG* or *NANOGP1* were distinct and well separated over the time course (Fig. 7G). Furthermore, by comparing the gene expression profiles with human embryo transcriptional data (Xiang et al., 2020), we further characterised the cell differentiation phenotype, and this also emphasised the differences after target gene depletion. *NANOG* knockdown naïve cells, starting from 4 days after doxycycline treatment, clustered with trophoblast and cytotrophoblast cells of the embryo, whereas at the earlier time-points (day 0 and day 2) *NANOG* knockdown naïve cells, and the non-induced cells and all *NANOGP1*-downregulated samples, instead clustered closer to pre- and early post-implantation epiblast (Fig. 7H). These data confirm that *NANOG* is required to maintain naïve pluripotency, and establish that *NANOG*-depleted naïve hPSCs have similar transcriptional profiles to trophoblast and cytotrophoblast lineages. In contrast to *NANOG*, the loss of *NANOGP1* expression does not disrupt the transcriptome of naïve pluripotent cells or cause

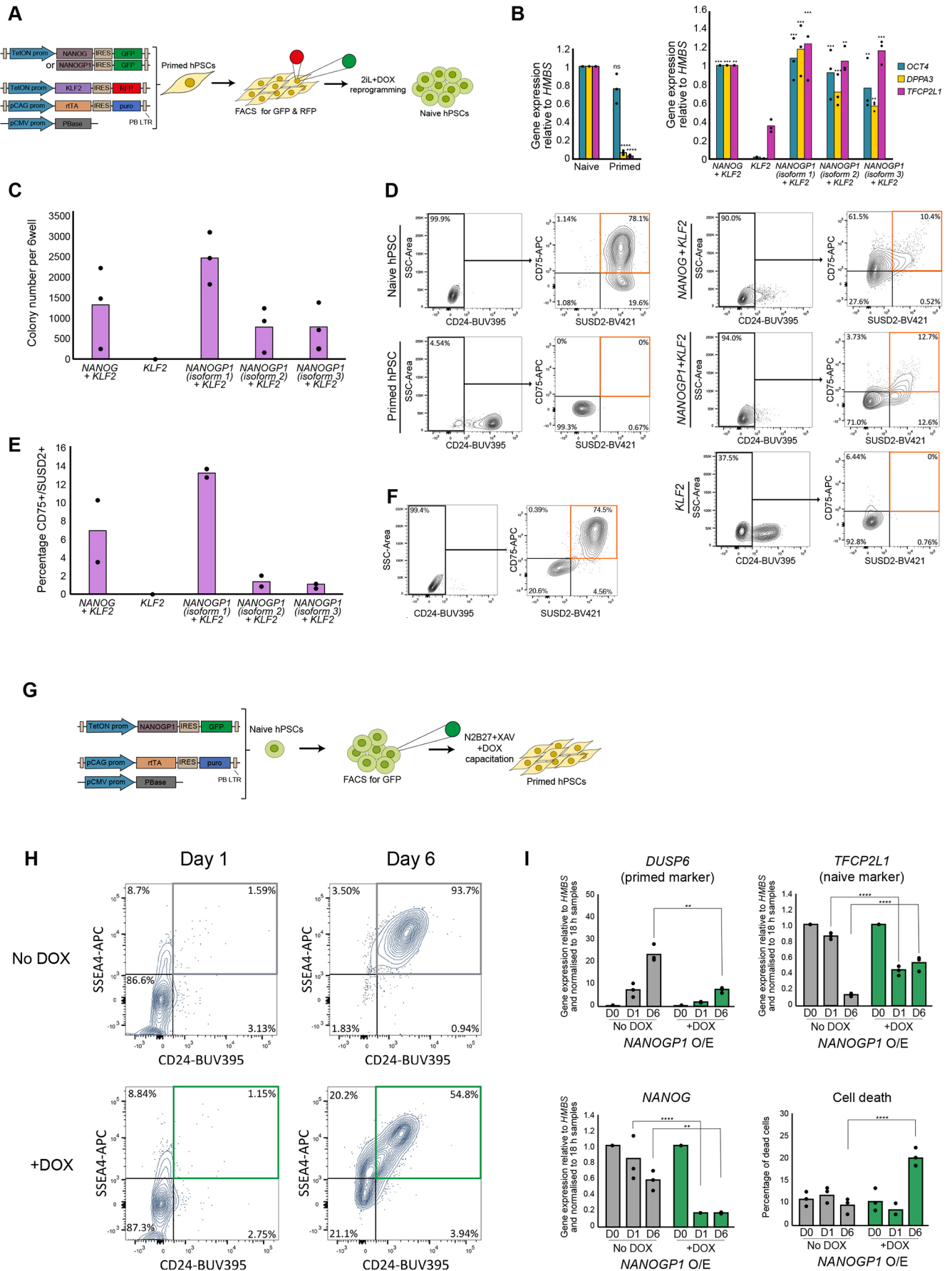


Fig. 6. See next page for legend.

Fig. 6. NANOGP1 is a strong inducer of naïve pluripotency.

(A) Experimental design for transgene-induced primed to naïve hPSC reprogramming. (B) Expression of pluripotency markers in established naïve and primed hPSCs (left), and in cultures after 12 days of DOX-induced reprogramming (right). RT-qPCR values are relative to *HMBS* expression and normalised to naïve hPSCs (left) and to the *NANOG+KLF2* sample (right). All three *NANOGP1* isoforms were tested. Mean and data points from three independent experiments are shown. Right: one-way ANOVA with Dunnett's multiple comparisons test compared all samples with the *KLF2*-only sample (* $P < 0.05$, ** $P < 0.005$, *** $P < 0.0005$, **** $P < 0.00005$). Left: unpaired *t*-test (two-tailed) compared the primed sample to the naïve samples (ns, not significant; **** $P < 0.00005$). (C) Number of alkaline phosphatase-positive colonies after 12 days of DOX-induced reprogramming. Mean and data points from three reprogramming experiments are shown. (D) Flow cytometry of cell-surface markers in established naïve and primed hPSCs, and in cultures after 12 days of DOX-induced reprogramming. Naïve hPSCs (CD24 negative; CD75 positive; SUSD2 positive) are in the upper right quadrant of the final gate. (E) Summary of flow cytometry data from D for two independent reprogramming experiments. (F) Stable cell-surface marker expression in established *NANOGP1+KLF2* (isoform 1) cell lines propagated in the absence of DOX in naïve hPSC medium for seven passages. (G) Experimental design for naïve to primed hPSC capacitation with enforced *NANOGP1* expression. (H) Cell-surface marker expression in cultures after 1 and 6 days of capacitation in the absence and presence of DOX. Primed hPSCs (CD24 positive; SSEA4 positive) are in the upper right quadrants. (I) Expression of marker genes in cultures at days 0, 1 and 6 of capacitation in the absence and presence of DOX. RT-qPCR values are relative to *HMBS* expression and normalised to day 0. Lower right: percentage of dead cells as measured using Trypan Blue staining. Mean and data points from three independent experiments are shown. An unpaired, two-tailed *t*-test compared the No DOX with +DOX samples at each timepoint (** $P < 0.005$, **** $P < 0.00005$; all other data are not significant).

trophectoderm differentiation. Additionally, *NANOGP1* does not provide functional redundancy for *NANOG*, as its expression was not sufficient to maintain naïve hPSCs in the absence of *NANOG*. In summary, these results demonstrate that downregulating the expression of *NANOG* in naïve hPSCs causes the loss of pluripotency, and that this function is not conserved for *NANOGP1*.

DISCUSSION

To better understand the role of pseudogenes in human development and pluripotency, we characterised and studied the function of *NANOGP1*, a tandem duplicate of the transcription factor *NANOG*. We found that *NANOGP1* has overlapping but distinct expression patterns with *NANOG* in stem cell states and human embryo development. The restricted expression profile in epiblast, germ cells and hPSCs prompted us to investigate whether *NANOGP1* could have conserved functional activities in naïve pluripotency. First, we found that *NANOGP1* has the capacity for gene auto-repression, as elevated expression of *NANOGP1* suppressed the expression of *NANOG* and *NANOGP1*. These findings additionally demonstrated that *NANOG* also has this function in human cells, which fulfils a prediction based on work in mouse pluripotent stem cells (Navarro et al., 2012). Second, *NANOGP1* was a strong inducer of naïve pluripotency when overexpressed in minimal reprogramming conditions, and was able to generate naïve hPSCs with comparable efficiency to *NANOG*. These results are consistent with the ability of *NANOG* orthologues, and moreover the *NANOG* homeodomain alone, to establish naïve pluripotency in mouse (Theunissen et al., 2011). The intact homeodomain of *NANOGP1*, and the presence of *NANOGP1* protein in human naïve pluripotent cells, therefore provide elevated levels of an active form of the key pluripotency factor *NANOG*. Notably, we found that the homeodomain sequence of *NANOGP1* has been disabled in other

primate species, further supporting the likelihood that this domain has been conserved in human and other great apes. Finally, because *NANOG* has dose-sensitive functions that are potentially mediated by concentration-dependent phase transitions (Choi et al., 2022), it is possible that *NANOGP1* might contribute to these effects by lowering the critical concentration that is required for *NANOG* to form condensates.

Despite these functional capabilities, we also found that *NANOGP1* is not required to maintain naïve pluripotency *in vitro*. By engineering cells that expressed gene-specific CRISPR-interference to transcriptionally repress *NANOGP1*, we found that naïve hPSCs were unaffected by the robust knockdown of *NANOGP1*. Interestingly, the capacity of *NANOGP1* to induce naïve pluripotency but be unnecessary for its maintenance parallels another naïve pluripotency factor – *KLF17* (Lea et al., 2021). In contrast, the knockdown of *NANOG* caused naïve hPSCs to exit the naïve state and differentiate towards the trophoblast lineage. This finding demonstrates that, unlike mouse naïve pluripotent stem cells (Chambers et al., 2007; Novo et al., 2016), human naïve cells require *NANOG*. It will be important to determine whether this requirement is related to the specific capacity of human naïve cells to differentiate into trophoblast (Castel et al., 2020; Cinkompumin et al., 2020; Dong et al., 2020; Guo et al., 2021; Io et al., 2021), which could underpin the different sensitivities to the loss of *NANOG*.

It is likely that the downregulation of *NANOGP1* has little effect in naïve hPSCs because *NANOG* remains robustly expressed. However, we cannot rule out subtle effects, including deficiencies after loss of *NANOGP1* that we have not yet identified. One interesting future direction would be to investigate whether the differences in predicted protein structures between *NANOGP1* and *NANOG* create functional or regulatory differences. A prominent difference between the predicted *NANOGP1* and *NANOG* proteins is a 39 amino acid deletion at the *NANOGP1* N terminus. The *NANOG* N terminus has a role in transcriptional interference by attracting co-repressors of cell differentiation, thereby opposing the transactivation role that is mediated by the C terminus (Chang et al., 2009). A key question, therefore, is whether *NANOGP1* might lack this co-repression activity. The *NANOG* N terminus is also a target for post-translational protein modifications, such as phosphorylation and ubiquitylation, and the control of protein turnover (Oh et al., 2005). Investigating the 39 amino acid deletion is particularly interesting from an evolutionary point of view. Both the N-terminal and C-terminal domains of mouse *NANOG* are involved in transcriptional transactivation (Chang et al., 2009; Do et al., 2009; Oh et al., 2005). In human, the N terminus loses this conserved function and remains less understood than its C-terminal counterpart. Therefore, studying molecular interactions of human *NANOG* protein has the potential to expand our understanding of pluripotency regulation in a human-specific context. Future studies could therefore be aimed at determining whether there are differences in protein stability and perdurance between *NANOG* and *NANOGP1*, and, by implication, whether *NANOGP1* might operate outside the processes that act to control and limit *NANOG* activity.

Previous predictions based on mutation analysis proposed that *NANOGP1* is ~22 million years old (Booth and Holland, 2004). Our comparative phylogenetic analysis of primate genome assemblies suggests an older duplication date, of either ~40 Mya, between the divergence of apes and Old World monkeys (25-35 Mya), and the earlier divergence of New World monkeys (40-50 Mya); or earlier, before the divergence of New World monkeys from other primates. The availability and in some cases the

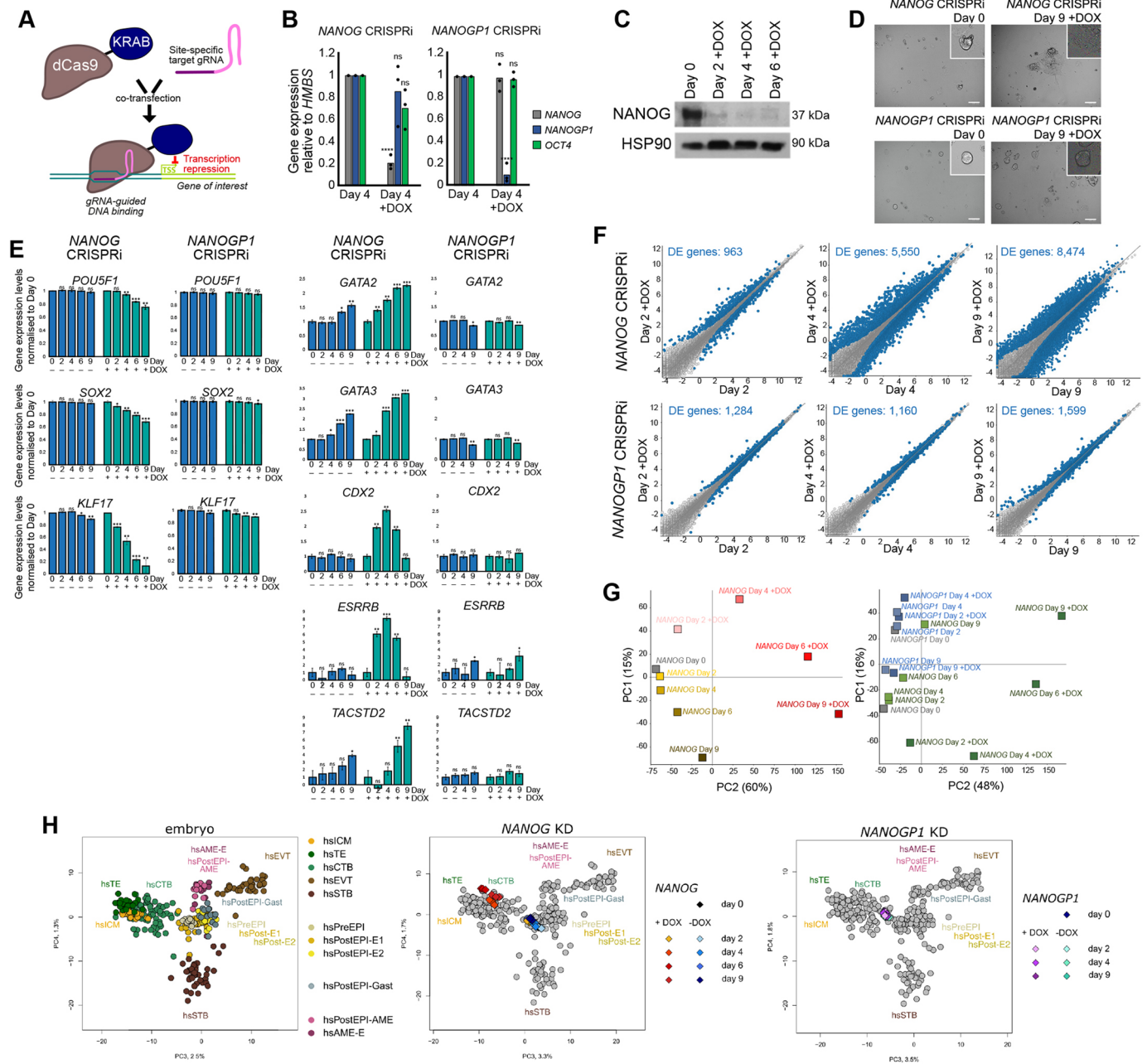


Fig. 7. NANOG is required to maintain naïve pluripotency, but NANOGP1 is dispensable. (A) DOX-inducible dCas9-KRAB CRISPRi to suppress *NANOG* and *NANOGP1* transcription in naïve hPSCs. (B) CRISPRi of *NANOG* (left) and *NANOGP1* (right) in naïve hPSCs. RT-qPCR values are relative to *HMBS* expression and normalised to day 4 samples. Mean and data points from three independent samples. An unpaired *t*-test (two-tailed) for each \pm DOX pair was performed (ns, not significant; *** $P < 0.00005$). (C) Reduced *NANOG* levels after DOX-induced *NANOG* CRISPRi in naïve hPSCs. (D) Bright-field images of *NANOG* and *NANOGP1* CRISPRi naïve hPSCs on day 0 and after 9 days of DOX treatment. Insets show representative colonies. Scale bars: 100 μ m. (E) Expression of undifferentiated (left) and trophectoderm markers (right) in *NANOG* and *NANOGP1* CRISPRi naïve hPSCs. Expression levels measured by RNA-seq are normalised to day 0 samples. Data are means \pm s.d. from three independent samples. An unpaired *t*-test (two-tailed) with multiple testing correction was performed between each timepoint and the corresponding day 0 sample (ns, not significant; * $P < 0.05$; ** $P < 0.005$; *** $P < 0.00005$). (F) Expression in *NANOG* (upper) and *NANOGP1* (lower) CRISPRi naïve hPSCs after DOX induction. Differentially expressed (DE) genes in blue [defined by a Wald test with Benjamini-Hochberg correction with a false discovery rate (FDR) of < 0.05]. (G) RNA-seq data of *NANOG* CRISPRi naïve hPSCs with and without DOX over a 9-day timecourse (left) and also with *NANOGP1* CRISPRi naïve hPSCs (right). Each data point is the average of three independent samples. (H) Left: transcriptionomes of annotated human embryo lineages (Xiang et al., 2020; Rostovskaya et al., 2022). On these maps, the transcriptomes of *NANOG* (centre) and *NANOGP1* (right) CRISPRi naïve hPSCs over a 9-day timecourse of DOX induction have been added. ICM, inner cell mass; TE, trophectoderm; CTB, cytotrophoblast; EVT, extravillous trophoblast; STB, syncytiotrophoblast; PreEPI, preimplantation epiblast; PostEPI, post-implantation epiblast; PostEPI-Gast, gastrulating stage; PostEPI-AME, post-implantation amniotic sac; AME, amniotic sac.

quality of current primate genome assemblies is insufficient to distinguish between the two scenarios, and this is a limitation of our study. More New World monkey and other primate genome

assemblies would be informative, and also it was not possible in most cases to search for the informative ‘scars’ that might remain after *NANOGP1* duplication and deletion. Therefore, it is only

possible at present to conclude that the duplication event took place at least ~40 Mya.

Our findings raise the question of why *NANOGPI* is retained in great apes but decayed in the genomes of lesser apes, Old World and New World monkeys? If *NANOGPI* provides epiblast cells with higher levels of NANOG-like activity, then perhaps this relates to, and is informative to understanding, the different developmental strategies between species. It is possible that the distinct modes of implantation (interstitial in great apes; superficial in New World and Old World monkeys), together with differences in the timing of blastocyst expansion and emergence of cell lineages, could point to a need to fine-tune transcription factor activities (Carter and Pijnenborg, 2011; Carter et al., 2015; Enders and Schlawke, 1986; Nakamura et al., 2016). To compare the functional role of transcription factors in early embryo development between different species, one future possibility could be the use of stem cell-derived embryo-like models (Kagawa et al., 2022; Liu et al., 2021; Sozen et al., 2021; Yanagida et al., 2021; Yu et al., 2021) from different species as a representative and genetically tractable system.

The majority of duplications in the human genome are segmental duplications, which, in particular, are thought to drive evolution of great apes and humans (Marques-Bonet et al., 2009a,b). *NANOGPI*, however, was formed by tandem duplication – an older evolutionarily mechanism. Strikingly, a tandem duplication of *NANOG* has occurred and was conserved at least twice: once, forming *NANOGPI*; and once, at a substantially earlier point, forming *NANOGNB*, which has diverged to such an extent that it was only recently recognised as a duplicate of *NANOG* (Dunwell and Holland, 2017). Independent *NANOG* duplications have also been reported in birds (Cañón et al., 2006), guinea pigs and some fish species (Scerbo et al., 2014). In all of these examples, the *NANOG* duplicates retain high similarity to their original ancestral sequences. These observations raise the possibility that the *NANOG*-containing region is somehow predisposed to duplication and retention of the duplication. In human, the chromosome region where *NANOG* is located also contains *DPPA3*, *POU5F1P3* and another pluripotency factor, *GDF3*, and collectively is called a ‘hotspot for teratocarcinoma’ owing to the high rate of chromosomal abnormalities (Clark et al., 2004; de Jong et al., 1990; Murty et al., 1990; Pain et al., 2005). Moreover, this region is also one of the most common amplification hotspots in hPSCs (International Stem Cell Initiative, 2011). There may be relevant parallels between the seemingly beneficial amplification of the *NANOG*-containing region throughout evolution and the aberrant amplification of the region associated with cell adaptation. A study in yeast showed that genes that are highly expressed before duplication have a higher chance of being retained for a longer evolutionary period and over a wider phylogenetic range (Mattenberger et al., 2017). If highly transcribed genes are more likely to be duplicated and retained, this raises specific and important implications for the genetic control of early epiblast and germ line development, particularly as chromosome changes in these cells would be heritable.

Pseudogenes are defined as disabled or defective versions of protein-coding genes and have long been considered as non-functional elements. The majority of pseudogenes in the human genome are processed. However, there are over 2000 unprocessed pseudogenes formed by duplication, many of which will have also copied their regulatory sequences. Careful annotation of pseudogenes, ideally supported by functional data, is important because they inform the reference list of genes and this impacts on whether sequence reads for the genes are mapped by default in

genome assemblies or are included in genetic screens and other related methods. Here, CRISPR-based approaches to epitope tag an endogenous pseudogene, and to recruit transcriptional repressive machinery to the endogenous promoter, enabled us to selectively explore pseudogene function. By doing this, we established that *NANOGPI* is protein coding and is expressed in pluripotent cells with functional activity. These results argue for the reclassification of *NANOGPI* as a protein-coding gene and for its reconsideration as a gene, rather than a pseudogene. In addition to *NANOGPI*, we found other highly expressed pseudogenes of prominent pluripotency factors, such as *POU5F1* and *DPPA3*, and it is therefore important to investigate whether they too are protein coding with functional properties. Defining pseudogene functionality and evolutionary conservation would help to uncover their involvement in species-specific developmental programmes and strategies.

MATERIALS AND METHODS

Human pluripotent stem cell lines

The use of human embryonic stem cells was carried out in accordance with approvals from the UK Stem Cell Bank Steering Committee. All cell lines used in this study were confirmed to be mycoplasma negative. Cell lines were not authenticated before use. WA09/H9 primed hPSCs were obtained from WiCell (Thomson et al., 1998). WA09/H9 NK2 (Takashima et al., 2014) and chemically-reset WA09/H9 (Guo et al., 2017) naïve hPSCs were kindly provided by Austin Smith (University of Exeter, UK). The CRISPRi Gen1B primed hPSCs (Mandegar et al., 2016) were kindly provided by Bruce Conklin and Li Gan (Gladstone Institutes, San Francisco, CA, USA).

Human pluripotent stem cell culture

All hPSC lines were maintained at 5% O₂ and 5% CO₂ at 37°C in a humidified incubator. Naïve hPSCs were cultured in N2B27 media composed of 1:1 DMEM/F12 and Neurobasal medium supplemented with 0.5× B-27, 0.5× N-2, 2 mM L-glutamine, 50 U/ml and 50 µg/ml penicillin-streptomycin and 0.1 mM β-mercaptoethanol (all ThermoFisher Scientific) and with 2 µM Gö6983 (Tocris), 1 µM PD0325901, 1 µM CHIR99021 and 20 ng/ml human LIF (all Wellcome-MRC Cambridge Stem Cell Institute) for t2iLGö medium (Takashima et al., 2014), or with 1 µM PD0325901, 2 µM Gö6983, 20 ng/ml human LIF and 2 µM XAV939 (Cell Guidance Systems) for PXGL medium (Bredenkamp et al., 2019b; Rostovskaya, 2022; Rostovskaya et al., 2019). Naïve hPSCs were grown either on irradiated MF1 mouse embryonic fibroblasts (MEFs) (Wellcome-MRC Cambridge Stem Cell Institute) on plates pre-coated with 0.1% gelatin (Sigma-Aldrich) or in feeder-free conditions using Geltrex Matrix (ThermoFisher Scientific) added to medium at a 1:300 dilution. Naïve hPSCs were passaged by 5 min incubation at 37°C with Accutase (BioLegend). Primed hPSCs were cultured on plates pre-treated with 5 µg/ml Vitronectin (ThermoFisher Scientific) in mTeSR Plus medium (STEMCELL Technologies) and passaged by 5 min incubation at room temperature with 0.5 mM EDTA in PBS.

NANOGPI epitope tagging

CRISPR/Cas12a-mediated gene editing, described previously (Zetsche et al., 2015), was adapted to epitope tag *NANOGPI*. Cas12a crRNA (IDT) targeting a region 10 bp upstream of the *NANOGPI* ATG site (5'-TGGCCCTGAAGAAAACCATCC-3') and a repair template containing an epitope tag (V5 or 3xFLAG; Table S2), were designed using CRISPOR (<http://crispor.tefor.net/>). For cell nucleofection, 5.6 µg Alt-R A.s. Cas12a crRNA and 40 µg Alt-R A.s. Cas12a Ultra protein were pre-assembled for 15 min at room temperature, combined with 2 µl 200 pmol/µl repair template (all reagents produced by IDT) and transfected into cR-H9 naïve hPSCs using a Neon Transfection System (ThermoFisher Scientific). Each transfection reaction was performed using 1 million cells per 100 µl Neon Transfection tip and with 1300 V, 30 ms and 1 pulse settings. After transfection, the cells were transferred to PXGL naïve hPSC media supplemented with 10 µM Y-27632 (Cell Guidance Systems). To improve

the rate of homology-directed repair, the cells were incubated in cold-shock conditions (32°C) for 24 h (Guo et al., 2018; Skarnes et al., 2019) at 5% O₂ and 5% CO₂ in a humidified incubator. Additionally, 2 μM M3814 (DNA-dependent protein kinase inhibitor) (Sigma-Aldrich) was added to the cell media for 72 h to repress non-homologous end-joining DNA repair (Riesenberg et al., 2019). To improve survival, 10 μM Y-27632 was added to the cells for 2 h before cell transfection and was kept in the media for 72 h after the transfection. The resultant cR-H9 NANOGP1-tag cell lines were expanded in PXGL media.

Inducible gene overexpression

To generate doxycycline-inducible gene overexpression vectors, gene cDNA was synthesised as a gBlocks Gene Fragment (IDT), cloned into a pCAG-IRES-Puro backbone vector (Niwa et al., 1991) and amplified with primers containing an *attB* sequence at their 5' ends (Table S3). The amplification product (*attB*-gene cDNA-*attB*) was cloned into a TetON-GFP/RFP plasmid kindly provided by Andras Nagy (Lunenfeld-Tanenbaum Research Institute, ON, Canada) (Woltjen et al., 2009) using a Gateway strategy (Hartley, 2003; Hartley et al., 2000) and was validated by Sanger sequencing (Genewiz). TetON plasmids, as well as plasmids encoding constitutively expressed reverse tetracycline-regulated transactivator gene (pCAG-rtTa-Puro) and a piggyBac transposase (pCyL43) (Wang et al., 2008) were transfected into primed H9 hPSCs using an Amaxa 4D nucleofactor (Lonza) with the setting CB-150. Stable cell lines were generated by 1 μg/ml puromycin selection for 48 h, followed by transient gene induction by adding 1 μM doxycycline for 48 h and flow sorting for fluorescent reporter expression. For all assays that included more than one cell line, the same sorting gate was used to sort reporter-positive cells in order to establish lines with similar gene expression levels.

Primed to naïve hPSC chemical reprogramming

Primed TetON-NANOGP1-GFP H9 hPSCs were reprogrammed into the naïve state using a chemical reprogramming method (Guo et al., 2017; Rugg-Gunn, 2022). Feeder-free cultures of primed hPSCs were passaged onto feeders in mTeSR Plus medium supplemented with 10 μM Y-27632 at a density of 10,000 per cm² (day 0) and provided with mTeSR Plus medium without Y-27632 on the following day. On day 2, the medium was changed to chemical reprogramming medium 1 (cRM-1), composed of N2B27 medium supplemented with 1 μM PD0325901, 10 ng/ml human LIF and 1 mM valproic acid sodium salt (Sigma-Aldrich). Starting from day 4, the medium was changed daily. On day 5, cRM-1 medium was replaced with chemical reprogramming medium 2 (cRM-2), composed of N2B27 medium supplemented with 1 μM PD0325901, 10 ng/ml human LIF, 2 μM Gö6983 and 2 μM XAV939. After several passages, the culture became homogeneous and was transferred to t2iLGö medium.

NANOGP1-mediated reprogramming

Primed H9 hPSC lines transfected with either *TetON-NANOGP1-GFP* (all three *NANOGP1* isoforms separately) plus *TetON-KLF2-RFP*, or with *TetON-NANOG-GFP* plus *TetON-KLF2-RFP*, were reprogrammed as described previously (Takashima et al., 2014). Before reprogramming, primed hPSCs were treated with 1 μM doxycycline for 24 h and flow-sorted for GFP⁺ signal or GFP⁺/RFP⁺ double-positive signal to establish transgenic lines with the equivalent level of reporter expression. Transgenic lines were then plated on feeders in KSR/FGF2 medium comprising 80% advanced DMEM, 20% knockout serum replacement (KSR), 2 mM L-glutamine, 50 U/ml and 50 μg/ml penicillin-streptomycin, 0.1 mM β-mercaptoethanol (all ThermoFisher Scientific) and 4 ng/ml basic fibroblast growth factor (Wellcome-MRC Cambridge Stem Cell Institute) supplemented with 10 μM Y-27632 (day 0) and, on the following day, the medium was changed to KSR/FGF2 supplemented with 1 μM doxycycline. On day 2, medium was changed to t2iL medium, composed of N2B27 medium with 1 μM PD0325901, 1 μM CHIR99021 and 10 ng/ml human LIF, supplemented with 1 μM doxycycline. t2iL medium was changed daily and cells were passaged every 5 days. On day 12, doxycycline was withdrawn and 5 μM Gö6983 was added. Reprogrammed cells were propagated in t2iLGö medium on feeders.

Naïve to primed hPSC capacitation

Naïve hPSCs were capacitated to a formative state as described by Rostovskaya et al. (2019). On day 0, naïve TetON-NANOGP1-1-GFP CR-H9 hPSCs were seeded in PXGL medium supplemented with 10 μM Y-27632 in feeder-free conditions on plates pre-coated with Geltrex at a seeding density of 16,000 per cm². On day 1, culture medium was replaced with PXGL without Y-27632. On day 2, medium was replaced with N2B27 supplemented with 2 μM XAV939, either with or without 1 μM doxycycline. Medium was then replaced every day and cells were passaged at a 1:2 ratio when 80% confluent. In total, hPSCs were cultured in N2B27 supplemented with XAV939 with or without doxycycline for 14 days.

Inducible gene expression knockdown

dCas9-iKRAB Gen1B CRISPRi *NANOGP1* and CRISPRi *NANOG* hPSC lines were generated as follows. Gene-specific gRNA oligonucleotides were phospho-annealed and cloned into pgRNA-CKB (pCAG-mKate2-T2A-bsd) vector (Mandegar et al., 2016), pre-digested with BsmBI (NEB) and pre-treated with FastAP (ThermoFisher Scientific). The *NANOGP1* gRNA sequence was designed and validated in this study, and the *NANOG* gRNA sequence was from Mandegar et al. (2016). Sequences are in Table S4. Linearised vector and phospho-annealed gRNA oligonucleotides were ligated at room temperature overnight with T4 DNA Ligase (ThermoFisher Scientific). Ligated products were validated by Sanger sequencing (Genewiz). Sequencing primers used were 5'-GAGATCCAGTTTGGTTAGTACCGGG-3' and 5'-ATGCATGGCGGTAATACGGTTAT-3'.

CRISPRi Gen1B primed hPSCs (Mandegar et al., 2016) were nucleofected with the *NANOGP1* and *NANOG* gRNA plasmids using Amaxa 4D Nucleofactor (setting CB-150), selected by blasticidin treatment (8 μg/ml for 5 days) and flow sorted for mKate2 expression. Primed CRISPRi Gen1B *NANOGP1* and *NANOG* lines were reprogrammed into the naïve state using 5i/L/A-mediated resetting (Fischer et al., 2022; Theunissen et al., 2014). To do this, primed feeder-free cultures were passaged onto feeders in mTeSR Plus medium supplemented with 10 μM Y-27632 at a density of 20,000 per cm² (day 0). On day 1, mTeSR Plus was replaced with 5i/L/A medium composed of N2B27 medium supplemented with 1 μM PD0325901, 20 ng/ml human LIF and 20 ng/ml activin A (Wellcome-MRC Cambridge Stem Cell Institute), 1 μM IM12, 0.5 μM SB590885, 10 μM Y-27632 and 1 μM WH-4-023 (all from Cell Guidance Systems). Cultures were passaged every 5 days and transferred to t2iLGö medium on day 18. CRISPRi was induced with 1 μM doxycycline.

Alkaline phosphatase activity

Colony formation assay was performed in combination with alkaline phosphatase (AP) staining (Štefková et al., 2015). Human PSCs were dissociated into single cells and plated into the experiment-specific medium onto feeders in six-well plates. On day 12, the cells were assayed for AP activity and imaged using a Zeiss Axio Observer Z1 with a 10× objective lens and Zeiss AxioVision software. Cells were fixed with 4% paraformaldehyde (PFA; Agar Scientific) in PBS, incubated in alkaline phosphatase staining solution (Merck) for 15 min and washed with PBS twice. The number of AP-positive colonies was counted.

Protein immunoprecipitation

All buffers used in this protocol were made with distilled water, were pre-chilled to 4°C and contained cOmplete EDTA-free protease inhibitor. All centrifugation steps were performed at 4°C. NANOGP1-V5 and NANOGP1-3xFLAG hPSCs were harvested and centrifuged for 5 min at 300 g, with 5 million cells per immunoprecipitation sample. To fractionate nuclei, pellets were resuspended in ice-cold buffer A [10 mM HEPES, 1.5 mM MgCl₂, 10 mM KCl, 0.5 mM DTT, 0.05% NP40 and 250 μ/ml benzonase nuclease (Sigma-Aldrich)], incubated for 10 min on ice and centrifuged for 10 min at 2000 g. Cell pellets were resuspended in 376 μl buffer B (5 mM HEPES, 1.5 mM MgCl₂, 0.2 mM EDTA, 0.5 mM DTT, 26% glycerol and 250 μ/ml benzonase nuclease), followed by 24 μl of 5 M NaCl. The resulting mix was homogenised using a Dounce on ice. Cell suspensions were kept on ice for 30 min followed by centrifugation for

20 min at 17,000 *g*. The supernatant was analysed by Bradford assay and stored on ice. Using a magnetic rack, protein A and protein G dynabeads (ThermoFisher Scientific) were washed twice with immunoprecipitation dilution buffer [150 mM Tris-HCl (pH 7.5), 150 mM NaCl and 0.5 mM EDTA]. Then, 5 µg of anti-V5 and anti-FLAG antibodies (Table S5) were added to the protein G and protein A magnetic beads, respectively, which were diluted in 500 µl immunoprecipitation dilution buffer. Tubes were kept on a rotating wheel at 4°C overnight. The next day, the beads were washed three times in the immunoprecipitation dilution buffer. Then, 475 µg (95%) of the nuclear protein obtained in the lysis step was added to the beads. 25 µg (5%) of each protein sample were set aside as input. Immunoprecipitation samples were rotated at 4°C overnight. The next day, beads were resuspended in the immunoprecipitation dilution buffer and washed for a total of three washes. To elute the immunoprecipitated complexes, beads were resuspended in 20 µl 5× protein loading dye and boiled at 75°C for 10 min. The eluate was diluted at 1× concentration, stored at -80°C and used in western blot assays.

Western blotting

Protein samples were extracted from frozen cell pellets, resuspended in ice-cold RIPA buffer (25 mM Tris/HCl, 140 mM NaCl, 1% Triton X-100, 0.5% SDS, 1 mM EDTA, 1 mM PMSF, 1 mM Na₃VO₄ and 1 mM NaF) supplemented with cOmplete protease inhibitor (Roche, 1836170). Cells were lysed by incubating on ice for 30 min. Lysates were centrifuged at 16,000 *g* for 30 min at 4°C. Protein concentration in supernatants was quantified using the Bradford assay. An appropriate volume of each lysate (containing 20–50 µg of the protein) was mixed with a 5× protein loading dye [5% β-mercaptoethanol, 0.02% bromophenol blue, 30% glycerol, 10% SDS and 250 mM Tris-Cl (pH 6.8)] and incubated at 90°C for 5 min. Samples were vortexed and placed on ice. Protein samples were run on a polyacrylamide vertical gel and transferred onto a polyvinylidene fluoride (PVDF) membrane using iBlot gel transfer system. The membrane was blocked with 5% milk (Sigma-Aldrich) in TBST (Tris-buffered saline+1% Tween 20) (Sigma-Aldrich) for 1 h at room temperature. Primary antibody was applied in TBST+5% milk overnight at 4°C. The next day, the membrane was washed three times with TBST and HRP-conjugated secondary antibody was applied for 1 h at room temperature. The membrane was washed three times and visualised by ECL or IRDye-conjugated secondary antibodies. Antibody details are provided in Table S5.

Immunofluorescence microscopy

Human PSCs were fixed in 12-well cell culture plates for 15 min at 4°C in 4% PFA in PBS, washed once with PBS and permeabilised with 0.4% Triton X-100 (Sigma-Aldrich) in PBS for 10 min at room temperature. Non-specific antibody binding was minimised by incubating cells with 3% BSA (Sigma-Aldrich)+0.1% Triton X-100/PBS for 1 h at room temperature. The cells were incubated with the appropriate primary antibody in 3% BSA+0.1% Triton X-100/PBS overnight at 4°C, before being washed four times with 0.1% Triton X-100/PBS and incubated with the appropriate secondary antibodies in 3% BSA+0.1% Triton X-100/PBS for 1 h at room temperature in the dark. Finally, the cells were washed three times in 0.1% Triton X-100/PBS [for nuclei staining, 1 µg/ml DAPI (Tocris) was added to the first wash] and twice in PBS. Wells were then filled with PBS, plates were sealed and stored at 4°C. Antibody details are provided in Table S6. Imaging was performed at the Babraham Institute Imaging Facility using a Nikon Live Cell Imager with a 20× objective lens.

Flow cytometry

Cells were dissociated with Accutase, washed with 2% FBS in PBS (wash buffer) and filtered through a 50 µm sterile strainer (Sysmex). Antibody labelling was performed by incubating cells in a Brilliant Stain Buffer (BD Biosciences) with antibodies for 30 min at 4°C in the dark. This was followed by a wash in wash buffer, cell pelleting at 300 *g* for 3 min and re-suspending the cells in 300 µl of the wash buffer. To identify live and dead cells, 0.1 µg/ml DAPI (Tocris) or Fixable Viability Dye eFluor 780 (eBioscience) was used. Antibody details are listed in Table S7. Flow cytometry analysis was performed on a BD LSR-Fortessa at the Babraham

Institute Flow Core. Cell-sorting experiments were performed on a BD Influx or a BD FACSAria Fusion. Data processing and downstream analysis were performed using FlowJo V10.1.

RNA sequencing

RNA was extracted using an RNeasy Mini Kit (Qiagen). Indexed libraries were made using 0.5 µg RNA per sample with NEBNext Ultra RNA Library Prep Kit for Illumina with the Poly(A) mRNA Magnetic Isolation Module (NEB) and NEBNext Multiplex Oligos for Illumina (NEB). An Agilent Bioanalyzer 2100 and KAPA Library Quantification Kit (KAPA Biosystems, KK4824) were used to identify library fragment size and concentration. Samples were sequenced as 75 bp single-end libraries on an Illumina NextSeq 500 at the Babraham Institute Genomics Facility, which generated 14–35 million uniquely mapped reads per library.

Sequencing files were analysed by FastQC v0.11.9 (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). RNA-sequencing reads were trimmed using Trim Galore v0.4.2 software (<https://github.com/FelixKrueger/TrimGalore>) to remove the adaptor sequences. Then, using HISAT2 v2.0.5 (Kim et al., 2019) guided by the Ensemble v70 gene models, trimmed reads were mapped to the human GRCh38 genome (Aken et al., 2016). Sequencing data were imported using Seqmonk software (<http://www.bioinformatics.babraham.ac.uk/projects/seqmonk/>). DESeq2 was used to identify genes expressed differentially (cut-off of $P < 0.05$ without independent filtering and after testing correction). To correct for the library size and variance among counts, regularised log transformation was applied before data visualisation. Principal component analysis (PCA) was performed using the top one thousand most variable genes across the experiment, and the 1st and 2nd PCs were plotted.

Polymerase chain reaction and genotyping primers

Polymerase chain reaction (PCR) was used to amplify various genomic and plasmid DNA fragments. PCR reactions were run in a BioRad Thermal Cycler T100. Polymerases Q5 HiFi (NEB), LongAmp Taq (NEB) and HotStarTaq (Qiagen) were used according to the manufacturer's instructions. Primer sequences used in PCR reactions, genotyping and DNA Sanger sequencing can be found in Table S8.

RT-qPCR

RNA was extracted using RNeasy Mini Kit (Qiagen) and then converted to cDNA using QuantiTect Reverse Transcription Kit (Qiagen). cDNA was diluted to 60 ng/µl and used in RT-qPCR using SYBR Green Jump Start Taq (Sigma-Aldrich) with 200 nM forward and reverse primers (Sigma-Aldrich; designed using Primer3 software (Untergasser et al., 2012)). Samples were run in technical triplicates in 96-well plates on a Bio-Rad CFX96 or in 384-well plates on a Bio-Rad CFX384. The results were analysed using the delta-delta cycle threshold method (relative quantity = $2^{-\Delta\Delta C_t}$) for which technical triplicates were averaged and normalised to the expression of a housekeeping gene *HMBS*. Data values represent mean ± s.d. of three biological replicates, unless stated otherwise. Statistical analyses are described in the figure legends. *NANOG* and *NANOGP1* expression in hPSCs was quantified using RT-qPCR primers, designed and validated to distinguish between the two genes. These two primer pairs, as well as other gene-specific primer sequences, can be found in Table S9.

Bioinformatics

Sequence comparison between pseudogenes and their ancestral genes

For each gene and pseudogene pair, the coding sequence of the gene and the transcript sequence of the pseudogene were extracted from the GRCh38 genome assembly based on the annotation in the Ensembl v108 annotation set. When the gene had multiple splice variants, the annotated Ensembl canonical transcript was used. The gene and pseudogene sequences were aligned using a global Needleman Wunsch alignment from the EMBOSS suite (v6.6.0) needle program (Madeira et al., 2022). Percentage identity was calculated between the first and last overlapping base pairs from the two sequences.

Identification of *NANOGP1* transcript variants

To identify putative *NANOGP1* transcripts, a combination of in-house-generated datasets of naïve hPSCs, as well as publicly available data from

Theunissen et al. (2016) (GEO accession number GSE84382), Pastor et al. (2016) (GEO accession number GSE76970) and Takashima et al. (2014) (ENA accession number PRJEB7132) was used. All raw data were processed with Trim Galore (Krueger et al., 2021) (adapter and quality trimming, v0.6.5) and mapped to the human GRCh38 genome using HISAT2 (v2.1.0; options `-dta -sp 1000,1000`), guided by known splice sites from Ensembl release 94 (Homo_sapiens.GRCh38.94.gtf).

To find evidence for splicing, aligned reads were first imported into SeqMonk (v1.43.1; Babraham Bioinformatics) as introns rather than exons, which effectively uses the CIGAR operation 'N' as the start and end coordinates of putative introns. Multi-mapping reads were filtered out (MAPQ \geq 20).

To identify likely exons, reads were then imported into SeqMonk as standard, i.e. spliced, RNA-seq reads (MAPQ \geq 20). Using read counts of exonic reads and introns identified as described above, the data were inspected and manually curated further to identify potential *NANOGPI* transcript variants. Transcript candidates appearing well supported by both exonic and intronic reads were termed *NANOGPI* isoforms 1-3 and taken forward for further analyses. GTF/GFF files were generated for *NANOGPI* isoforms 1-3 and were included as additional annotations for both HISAT2 mapping and further analyses in SeqMonk.

To identify potential open reading frames of *NANOGPI* isoforms 1-3, their hypothetical cDNA sequences were then screened for open reading frames (ORF) using the NCBI Open Reading Frame Finder tool (<https://www.ncbi.nlm.nih.gov/orffinder/>). The longest ORFs, resulting in predicted proteins between 255 and 266 amino acids in length, were taken forward for multiple sequence alignments (Madeira et al., 2022) and additional analyses.

Disambiguation of *NANOG* and *NANOGPI*

To investigate the cross-mapping of reads from the *NANOG* to the *NANOGPI* locus, and vice versa, cDNA sequences for *NANOG* (NANOG-201, Ensembl) and *NANOGPI* (isoform 1) were used and converted to simulated FastQ files [as 43 bp (as in Petropoulos et al., 2016) or 100 bp single-end reads, in steps of 1 bp from start to end]. These *NANOG* and *NANOGPI* FastQ files were then aligned to the human GRCh38 genome (using HISAT2, v2.1.0; Kim et al., 2019); the amount of cross-mapping was either negligible or non-existent for unfiltered or multi-mapping filtered (MAPQ \geq 20) reads, respectively.

Human embryo data processing

The RNA-seq data of 1481 human embryo single cells from Petropoulos et al. (2016) were downloaded (accession number ERP012552) and categorised into the following groups: 8c, MOR, eICM, eTE, EPI, TE, PE, eUndef, Inter. Cell annotations were taken from Stirparo et al. (2018). The data were mapped to the human GRCh38 genome using HISAT2 (v2.1.0; Kim et al., 2019) using options `-dta -sp 1000,1000`, guided by known splice sites from Ensembl release 94 (Homo_sapiens.GRCh38.94.gtf) to which a custom *NANOGPI* mRNA annotation had been added manually. Reads were then filtered for unique alignments (MAPQ $>$ 20), and log₂ RPM counts for genes were calculated with SeqMonk (v1.43.1; Babraham Bioinformatics; assuming non-strand specific libraries and merging transcript isoforms). Violin plots of expression values for genes of interest were then calculated for different developmental stages using the beanplot library and RStudio (v1.1.463).

The RNA-seq data of 557 human embryo single cells from Xiang et al. (2020) were downloaded (accession number GSE136447) and categorised into the following groups: ICM, EPI, PrE and TrB. The data were mapped to the human GRCh38 genome using HISAT2 (v2.1.0; Kim et al., 2019) using options `-dta -sp 1000,1000`, guided by known splice sites from Ensembl release 94 (Homo_sapiens.GRCh38.94.gtf) to which a custom *NANOGPI* mRNA annotation had been added manually. Reads were then filtered for unique alignments (MAPQ $>$ 20) and log₂ RPM counts for genes were calculated with SeqMonk (v1.43.1; Babraham Bioinformatics; assuming non-strand specific libraries and merging transcript isoforms). Violin plots of expression values for genes of interest were then calculated for different epiblast developmental stages in R (RStudio).

Evolutionary genetics

To investigate the genomic structure of the *NANOG/NANOGPI* locus throughout evolution, the most recent assemblies of nine primate species (Table S10) were analysed. Approximate genomic coordinates of *NANOG* and *NANOGPI* (if present) were identified using BLAST (basic local alignment search tool; Sayers et al., 2022) and Needle (Madeira et al., 2022) pairwise sequence alignment tools. Within each assembly, a ~250 kb genomic region, including *NANOG*, *NANOGPI* and their surrounding genes was extracted. The *NANOGPI* open reading frame for each species was also extracted. DNA and its corresponding amino acid sequences of *NANOG* and *NANOGPI* were aligned using MEGA (Tamura et al., 2007) and ClustalW (Madeira et al., 2022). Codeml and codonml PAML (v4.8a) programs were run for the phylogenetic analysis of amino acid sequences with maximum likelihood under M0, M1, M7 and M8 models (Yang and Nielsen, 2000). Dotter (Barson and Griffiths, 2016) and Miropeats (Parsons, 1995) were used for visualising the *NANOG/NANOGPI* duplication site, detecting boundaries of the duplicated region and measuring conservation/divergence between the duplicated sequences since the duplication event.

The Gibbon nomLeu3.0 assembly was found to be not suitable for investigating the *NANOG* region due to having large gaps in the relevant region. To resolve this, unpublished gibbon genome assembly data based on long-read sequencing, kindly provided by Evan Eichler (University of Washington), was analysed. To visualise the *NANOG*-containing locus, human *NANOG* and *NANOGPI* sequence was mapped to gibbon contigs using Minimap2 (Li, 2018; Parsons, 1995).

For GC content calculation, enhancer regions were first extracted from human genome assembly (GRCh38 build) as FASTA files based on previously provided genomic coordinates. We then calculated GC content by dividing the sum of G and C nucleotide counts (G+C) to the total nucleotide count (G+C+T+A) at a genomic region. We used a 30 base-pair sliding-window approach to calculate GC content along the enhancer regions, and plotted GC percentages against genomic coordinates.

Statistics and reproducibility

Sample size was not predefined. Samples were randomly allocated to experimental groups by the investigator. All experiments were replicated at least three times using independent biological samples. All images are representative. Data points were collected without investigator blinding. No data were excluded. Graphs were prepared using R and Prism v8. *P*-values were calculated as specified in figure legends.

Acknowledgements

We thank Austin Smith, Bruce Conklin and Li Gan for providing cell lines, and Evan Eichler for providing access to unpublished gibbon genome data. We are grateful to Paula Kokko-Gonzales and Amelia Edwards at the Babraham Institute Genomics Facility; to Rachael Walker, Rebecca Roberts and the team at the Babraham Institute Flow Core; to Simon Andrews from the Babraham Bioinformatics Group; and to the Wellcome – MRC Cambridge Stem Cell Institute Tissue Culture Facility for providing reagents.

Competing interests

The authors declare no competing or financial interests.

Author contributions

Conceptualization: K.M., P.J.R.-G.; Methodology: A.N.; Formal analysis: K.M., G.A., F.K., J.W., M.R., C.K., P.J.R.-G.; Investigation: K.M., G.A., J.W., M.R., A.B., S.W., P.J.R.-G.; Data curation: F.K.; Writing - original draft: K.M., G.A., P.J.R.-G.; Writing - review & editing: K.M., G.A., F.K., J.W., M.R., A.N., A.B., C.K., S.W., A.S., P.J.R.-G.; Visualization: K.M., G.A., J.W., M.R., P.J.R.-G.; Supervision: A.S., P.J.R.-G.; Project administration: P.J.R.-G.; Funding acquisition: A.S., P.J.R.-G.

Funding

This research was supported by the Biotechnology and Biological Sciences Research Council (BBS/E/B/000C0421, BBS/E/B/000C0422 to P.J.R.-G.; Core Capability Grant), the Medical Research Council (MR/T011769/1 and MR/V02969X/1 to P.J.R.-G.), the Wellcome Trust (215116/Z/18/Z to P.J.R.-G.), the Darwin Trust (K.M.), the Cambridge Commonwealth, European and International Trust (K.M.), the Cambridge Biosciences Biotechnology and Biological Sciences Research Council DTP (to A.B.) and Erasmus+ (EU programme for education, training, youth and sport

to G.A. and A.S.). Open Access funding provided by the Babraham Institute. Deposited in PMC for immediate release.

Data availability

RNA sequencing datasets have been deposited in the GEO under the accession number GSE204934.

Peer review history

The peer review history is available online at <https://journals.biologists.com/dev/lookup/doi/10.1242/dev.201155.reviewer-comments.pdf>

References

- Aken, B. L., Achuthan, P., Akanni, W., Amode, M. R., Bernsdorff, F., Bhai, J., Billis, K., Carvalho-Silva, D., Cummins, C., Clapham, P. et al. (2016). Ensembl 2017. *Nucleic Acids Res.* **45**, D635–D642. doi:10.1093/nar/gkw1104
- Bailey, J. A., Gu, Z., Clark, R. A., Reinert, K., Samonte, R. V., Schwartz, S., Adams, M. D., Myers, E. W., Li, P. W. and Eichler, E. E. (2002). Recent segmental duplications in the human genome. *Science* **297**, 1003–1007. doi:10.1126/science.1072047
- Barson, G. and Griffiths, E. (2016). SeqTools: visual tools for manual analysis of sequence alignments. *BMC Res. Notes* **9**, 39. doi:10.1186/s13104-016-1847-3
- Booth, H. A. F. and Holland, P. W. H. (2004). Eleven daughters of NANOG*. *Genomics* **84**, 229–238. doi:10.1016/j.ygeno.2004.02.014
- Bredenkamp, N., Stirparo, G. G., Nichols, J., Smith, A. and Guo, G. (2019a). The cell-surface marker sushi containing domain 2 facilitates establishment of human naive pluripotent stem cells. *Stem Cell Rep.* **12**, 1212–1222. doi:10.1016/j.stemcr.2019.03.014
- Bredenkamp, N., Yang, J., Clarke, J., Stirparo, G. G., von Meyenn, F., Dietmann, S., Baker, D., Drummond, R., Ren, Y., Li, D. et al. (2019b). Wnt inhibition facilitates RNA-mediated reprogramming of human somatic cells to naive pluripotency. *Stem Cell Rep.* **13**, 1083–1098. doi:10.1016/j.stemcr.2019.10.009
- Cañón, S., Herranz, C. and Manzanares, M. (2006). Germ cell restricted expression of chick Nanog. *Dev. Dyn.* **235**, 2889–2894. doi:10.1002/dvdy.20927
- Carter, A. M. and Pijnborg, R. (2011). Evolution of invasive placentation with special reference to non-human primates. *Best Pract. Res. Clin. Obstet. Gynaecol.* **25**, 249–257. doi:10.1016/j.bpobgyn.2010.10.010
- Carter, A. M., Enders, A. C. and Pijnborg, R. (2015). The role of invasive trophoblast in implantation and placentation of primates. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **370**, 20140070. doi:10.1098/rstb.2014.0070
- Castel, G., Meistermann, D., Bretin, B., Firmin, J., Blin, J., Loubersac, S., Bruneau, A., Chevolleau, S., Kilens, S., Chariou, C. et al. (2020). Induction of human trophoblast stem cells from somatic cells and pluripotent stem cells. *Cell Rep.* **33**, 108419. doi:10.1016/j.celrep.2020.108419
- Chambers, I., Colby, D., Robertson, M., Nichols, J., Lee, S., Tweedie, S. and Smith, A. (2003). Functional expression cloning of Nanog, a pluripotency sustaining factor in embryonic stem cells. *Cell* **113**, 643–655. doi:10.1016/S0092-8674(03)00392-1
- Chambers, I., Silva, J., Colby, D., Nichols, J., Nijmeijer, B., Robertson, M., Vrana, J., Jones, K., Grotewold, L. and Smith, A. (2007). Nanog safeguards pluripotency and mediates germline development. *Nature* **450**, 1230–1234. doi:10.1038/nature06403
- Chang, C. P., Brocchieri, L., Shen, W. F., Largman, C. and Cleary, M. L. (1996). Pbx modulation of Hox homeodomain amino-terminal arms establishes different DNA-binding specificities across the Hox locus. *Mol. Cell. Biol.* **16**, 1734–1745. doi:10.1128/MCB.16.4.1734
- Chang, D. F., Tsai, S. C., Wang, X. C., Xia, P., Senadheera, D. and Lutzko, C. (2009). Molecular characterization of the human NANOG protein. *Stem Cells* **27**, 812–821. doi:10.1634/stemcells.2008-0657
- Charrier, C., Joshi, K., Coutinho-Budd, J., Kim, J.-E., Lambert, N., de Marchena, J., Jin, W.-L., Vanderhaeghen, P., Ghosh, A., Sassa, T. et al. (2012). Inhibition of SRGAP2 function by its human-specific paralogs induces neoteny during spine maturation. *Cell* **149**, 923–935. doi:10.1016/j.cell.2012.03.034
- Choi, K.-J., Quan, M. D., Qi, C., Lee, J.-H., Tsoi, P. S., Zahabiyon, M., Bajic, A., Hu, L., Prasad, B. V. V., Liao, S.-C. J. et al. (2022). NANOG prion-like assembly mediates DNA bridging to facilitate chromatin reorganization and activation of pluripotency. *Nat. Cell Biol.* **24**, 737–747. doi:10.1038/s41556-022-00896-x
- Chovanec, P., Collier, A. J., Krueger, C., Várnai, C., Semplich, C. I., Schoenfelder, S., Corcoran, A. E. and Rugg-Gunn, P. J. (2021). Widespread reorganisation of pluripotent factor binding and gene regulatory interactions between human pluripotent states. *Nat. Commun.* **12**, 2098. doi:10.1038/s41467-021-22201-4
- Cinkornpumin, J. K., Kwon, S. Y., Guo, Y., Hossain, I., Sirois, J., Russett, C. S., Tseng, H.-W., Okae, H., Arima, T., Duchaine, T. F. et al. (2020). Naive human embryonic stem cells can give rise to cells with a trophoblast-like transcriptome and methylome. *Stem Cell Rep.* **15**, 198–213. doi:10.1016/j.stemcr.2020.06.003
- Clark, A. T., Rodriguez, R. T., Bodnar, M. S., Abeyta, M. J., Cedars, M. I., Turek, P. J., Firpo, M. T. and Reijo Pera, R. A. (2004). Human *STELLAR*, *NANOG*, and *GDF3* Genes Are Expressed in Pluripotent Cells and Map to Chromosome 12p13, a Hotspot for Teratocarcinoma. *Stem Cells* **22**, 169–179. doi:10.1634/stemcells.22-2-169
- Collier, A. J., Panula, S. P., Schell, J. P., Chovanec, P., Reyes, A. P., Petropoulos, S., Corcoran, A. E., Walker, R., Douagi, I., Lanner, F. et al. (2017). Comprehensive cell surface protein profiling identifies specific markers of human naive and primed pluripotent states. *Cell Stem Cell* **20**, 874–890.e7. doi:10.1016/j.stem.2017.02.014
- de Jong, B., Wolter Oosterhuis, J., Castedo, S. M. M. J., Vos, A. and te Meerman, G. J. (1990). Pathogenesis of adult testicular germ cell tumors. *Cancer Genet. Cytogenet.* **48**, 143–167. doi:10.1016/0165-4608(90)90115-Q
- Dennis, M. Y. and Eichler, E. E. (2016). Human adaptation and evolution by segmental duplication. *Curr. Opin. Genet. Dev.* **41**, 44–52. doi:10.1016/j.gde.2016.08.001
- Do, H.-J., Lee, W.-Y., Lim, H. Y., Oh, J.-H., Kim, D.-K., Kim, J.-H., Kim, T. and Kim, J.-H. (2009). Two potent transactivation domains in the C-terminal region of human NANOG mediate transcriptional activation in human embryonic carcinoma cells. *J. Cell. Biochem.* **106**, 1079–1089. doi:10.1002/jcb.22089
- Dong, C., Beltcheva, M., Gontarz, P., Zhang, B., Popli, P., Fischer, L. A., Khan, S. A., Park, K.-M., Yoon, E.-J., Xing, X. et al. (2020). Derivation of trophoblast stem cells from naive human pluripotent stem cells. *eLife* **9**, e52504. doi:10.7554/eLife.52504
- Dunwell, T. L. and Holland, P. W. H. (2017). A sister of *NANOG* regulates genes expressed in pre-implantation human development. *Open Biol.* **7**, 170027. doi:10.1098/rsob.170027
- Eberle, I., Pless, B., Braun, M., Dingermann, T. and Marschalek, R. (2010). Transcriptional properties of human NANOG1 and NANOG2 in acute leukemic cells. *Nucleic Acids Res.* **38**, 5384–5395. doi:10.1093/nar/gkq307
- Enders, A. C. and Schlafke, S. (1986). Implantation in nonhuman primates and in the human. In *Comparative Primate Biology, vol. 3: Reproduction and Development* (ed. W. R. Dukelow and J. Erwin), pp. 291–310. New York, NY: Alan R. Liss Inc.
- Fairbanks, D. J. and Maughan, P. J. (2006). Evolution of the NANOG pseudogene family in the human and chimpanzee genomes. *BMC Evol. Biol.* **6**, 12. doi:10.1186/1471-2148-6-12
- Fares, M. A. (2014). The evolution of protein moonlighting: adaptive traps and promiscuity in the chaperonins. *Biochem. Soc. Trans.* **42**, 1709–1714. doi:10.1042/BST20140225
- Fischer, L. A., Khan, S. A. and Theunissen, T. W. (2022). Induction of human naive pluripotency using 5i/L/A medium. *Methods Mol. Biol.* **2416**, 13–28. doi:10.1007/978-1-0716-1908-7_2
- Florio, M., Albert, M., Taverna, E., Namba, T., Brandl, H., Lewitus, E., Haffner, C., Sykes, A., Wong, F. K., Peters, J. et al. (2015). Human-specific gene ARHGAP11B promotes basal progenitor amplification and neocortex expansion. *Science* **347**, 1465–1470. doi:10.1126/science.aaa1975
- Force, A., Lynch, M., Bryan Pickett, F., Amores, A., Yan, Y.-L. and Postlethwait, J. (1999). Preservation of duplicate genes by complementary, degenerative mutations. *Genetics* **151**, 1531–1545. doi:10.1093/genetics/151.4.1531
- Gkoutela, S., Zhang, K. X., Shafiq, T. A., Liao, W.-W., Hargan-Calvopiña, J., Chen, P.-Y. and Clark, A. T. (2015). DNA demethylation dynamics in the human prenatal germline. *Cell* **161**, 1425–1436. doi:10.1016/j.cell.2015.05.012
- Guo, G., von Meyenn, F., Santos, F., Chen, Y., Reik, W., Bertone, P., Smith, A. and Nichols, J. (2016). Naive pluripotent stem cells derived directly from isolated cells of the human inner cell mass. *Stem Cell Rep.* **6**, 437–446. doi:10.1016/j.stemcr.2016.02.005
- Guo, G., von Meyenn, F., Rostovskaya, M., Clarke, J., Dietmann, S., Baker, D., Sahakyan, A., Myers, S., Bertone, P., Reik, W. et al. (2017). Epigenetic resetting of human pluripotency. *Development* **144**, 2748–2763. doi:10.1242/dev.146811
- Guo, Q., Mintier, G., Ma-Edmonds, M., Storton, D., Wang, X., Xiao, X., Kienzle, B., Zhao, D. and Feder, J. N. (2018). “Cold shock” increases the frequency of homology directed repair gene editing in induced pluripotent stem cells. *Sci. Rep.* **8**, 2080. doi:10.1038/s41598-018-20358-5
- Guo, G., Stirparo, G. G., Strawbridge, S. E., Spindlow, D., Yang, J., Clarke, J., Dattani, A., Yanagida, A., Li, M. A., Myers, S. et al. (2021). Human naive epiblast cells possess unrestricted lineage potential. *Cell Stem Cell* **28**, 1040–1056.e6. doi:10.1016/j.stem.2021.02.025
- Hart, A. H., Hartley, L., Ibrahim, M. and Robb, L. (2004). Identification, cloning and expression analysis of the pluripotency promoting Nanog genes in mouse and human. *Dev. Dyn.* **230**, 187–198. doi:10.1002/dvdy.20034
- Hartley, J. L. (2003). Use of the gateway system for protein expression in multiple hosts. *Curr. Protoc. Protein Sci.* doi:10.1002/0471140864.ps0517s30
- Hartley, J. L., Temple, G. F. and Brasch, M. A. (2000). DNA cloning using in vitro site-specific recombination. *Genome Res.* **10**, 1788–1795. doi:10.1101/gr.143000
- Hyslop, L., Stojkovic, M., Armstrong, L., Walter, T., Stojkovic, P., Przyborski, S., Herbert, M., Murdoch, A., Strachan, T. and Lako, M. (2005). Downregulation of NANOG Induces Differentiation of Human Embryonic Stem Cells to

- Extraembryonic Lineages. *Stem Cells* **23**, 1035-1043. doi:10.1634/stemcells.2005-0080
- International Stem Cell Initiative.** (2011). Screening ethnically diverse human embryonic stem cells identifies a chromosome 20 minimal amplicon conferring growth advantage. *Nat. Biotechnol.* **29**, 1132-1144. doi:10.1038/nbt.2051
- Io, S., Kabata, M., Iemura, Y., Semi, K., Morone, N., Minagawa, A., Wang, B., Okamoto, I., Nakamura, T., Kojima, Y. et al.** (2021). Capturing human trophoblast development with naive pluripotent stem cells in vitro. *Cell Stem Cell* **28**, 1023-1039.e13. doi:10.1016/j.stem.2021.03.013
- Kagawa, H., Javali, A., Khoei, H. H., Sommer, T. M., Sestini, G., Novatchkova, M., Scholte Op Reimer, Y., Castel, G., Bruneau, A., Maenhoudt, N. et al.** (2022). Human blastoids model blastocyst development and implantation. *Nature* **601**, 600-605. doi:10.1038/s41586-021-04267-8
- Kim, D., Paggi, J. M., Park, C., Bennett, C. and Salzberg, S. L.** (2019). Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**, 907-915. doi:10.1038/s41587-019-0201-4
- Kondrashov, F. A. and Kondrashov, A. S.** (2006). Role of selection in fixation of gene duplications. *J. Theor. Biol.* **239**, 141-151. doi:10.1016/j.jtbi.2005.08.033
- Krueger, F., James, F., Ewels, P., Afyounian, E. and Schuster-Boeckler, B.** (2021). TrimGalore: v0.6.7 - Zenodo. <https://doi.org/10.5281/zenodo.5127899>
- Lea, R. A., McCarthy, A., Boeing, S., Fallesen, T., Elder, K., Snell, P., Christie, L., Adkins, S., Shaikly, V., Taranissi, M. et al.** (2021). KLF17 promotes human naive pluripotency but is not required for its establishment. *Development* **148**, dev199378. doi:10.1242/dev.199378
- Li, H.** (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094-3100. doi:10.1093/bioinformatics/bty191
- Lie, K.-H., Tuch, B. E. and Sidhu, K. S.** (2012). Suppression of NANOG induces efficient differentiation of human embryonic stem cells to pancreatic endoderm. *Pancreas* **41**, 54-64. doi:10.1097/MPA.0b013e31822362e4
- Liu, X., Tan, J. P., Schröder, J., Aberkane, A., Ouyang, J. F., Mohenska, M., Lim, S. M., Sun, Y. B. Y., Chen, J., Sun, G. et al.** (2021). Modelling human blastocysts by reprogramming fibroblasts into iBlastoids. *Nature* **591**, 627-632. doi:10.1038/s41586-021-03372-y
- Madeira, F., Pearce, M., Tivey, A. R. N., Basutkar, P., Lee, J., Edbali, O., Madhusoodanana, N., Kolesnikov, A. and Lopez, R.** (2022). Search and sequence analysis tools services from EMBL-EBI in 2022. *Nucleic Acids Res.* **50**, W276-W279. doi:10.1093/nar/gkac240
- Magadum, S., Banerjee, U., Murugan, P., Gangapur, D. and Ravikesavan, R.** (2013). Gene duplication as a major force in evolution. *J. Genet.* **92**, 155-161. doi:10.1007/s12041-013-0212-8
- Mandegar, M. A., Huebsch, N., Frolov, E. B., Shin, E., Truong, A., Olvera, M. P., Chan, A. H., Miyaoka, Y., Holmes, K., Ian Spencer, C. et al.** (2016). CRISPR interference efficiently induces specific and reversible gene silencing in human iPSCs. *Cell Stem Cell* **18**, 541-553. doi:10.1016/j.stem.2016.01.022
- Marques-Bonet, T., Girirajan, S. and Eichler, E. E.** (2009a). The origins and impact of primate segmental duplications. *Trends Genet.* **25**, 443-454. doi:10.1016/j.tig.2009.08.002
- Marques-Bonet, T., Kidd, J. M., Ventura, M., Graves, T. A., Cheng, Z., Hillier, L. D. W., Jiang, Z., Baker, C., Malfavon-Borja, R., Fulton, L. A. et al.** (2009b). A burst of segmental duplications in the genome of the African great ape ancestor. *Nature* **457**, 877-881. doi:10.1038/nature07744
- Mattenberger, F., Sabater-Muñoz, B., Toft, C. and Fares, M. A.** (2017). The phenotypic plasticity of duplicated genes in *Saccharomyces cerevisiae* and the origin of adaptations. *G3* **7**, 63-75. doi:10.1534/g3.116.035329
- Mullin, N. P., Varghese, J., Colby, D., Richardson, J. M., Findlay, G. M. and Chambers, I.** (2021). Phosphorylation of NANOG by casein kinase I regulates embryonic stem cell self-renewal. *FEBS Lett.* **595**, 14-25. doi:10.1002/1873-3468.13969
- Murty, V. V. S., Dmitrovsky, E., Bosl, G. J. and Chaganti, R. S. K.** (1990). Nonrandom chromosome abnormalities in testicular and ovarian germ cell tumor cell lines. *Cancer Genet. Cytogenet.* **50**, 67-73. doi:10.1016/0165-4608(90)90239-7
- Nakamura, T., Okamoto, I., Sasaki, K., Yabuta, Y., Iwatani, C., Tsuchiya, H., Seita, Y., Nakamura, S., Yamamoto, T. and Saitou, M.** (2016). A developmental coordinate of pluripotency among mice, monkeys and humans. *Nature* **537**, 57-62. doi:10.1038/nature19096
- Navarro, P., Festuccia, N., Colby, D., Gagliardi, A., Mullin, N. P., Zhang, W., Karwacki-Neisius, V., Osorno, R., Kelly, D., Robertson, M. et al.** (2012). OCT4/SOX2-independent *Nanog* autorepression modulates heterogeneous *Nanog* gene expression in mouse ES cells. *EMBO J.* **31**, 4547-4562. doi:10.1038/emboj.2012.321
- Niwa, H., Yamamura, K.-I. and Miyazaki, J.-I.** (1991). Efficient selection for high-expression transfectants with a novel eukaryotic vector. *Gene* **108**, 193-199. doi:10.1016/0378-1119(91)90434-D
- Novo, C. L., Tang, C., Ahmed, K., Djuric, U., Fussner, E., Mullin, N. P., Morgan, N. P., Hayre, J., Sienerth, A. R., Elderkin, S. et al.** (2016). The pluripotency factor Nanog regulates pericentromeric heterochromatin organization in mouse embryonic stem cells. *Genes Dev.* **30**, 1101-1115. doi:10.1101/gad.275685.115
- Oh, J.-H., Do, H.-J., Yang, H.-M., Moon, S.-Y., Cha, K.-Y., Chung, H.-M. and Kim, J.-H.** (2005). Identification of a putative transactivation domain in human Nanog. *Exp. Mol. Med.* **37**, 250-254. doi:10.1038/emmm.2005.33
- Ohta, T.** (2000). Evolution of gene families. *Gene* **259**, 45-52. doi:10.1016/S0378-1119(00)00428-5
- Pain, D., Chirn, G.-W., Strassel, C. and Kemp, D. M.** (2005). Multiple Retroseudogenes from Pluripotent Cell-specific Gene Expression Indicates a Potential Signature for Novel Gene Identification. *J. Biol. Chem.* **280**, 6265-6268. doi:10.1074/jbc.C400587200
- Parsons, J. D.** (1995). Miropeats: graphical DNA sequence comparisons. *Comput. Appl. Biosci.* **11**, 615-619. doi:10.1093/bioinformatics/11.6.615
- Pastor, W. A., Chen, D., Liu, W., Kim, R., Sahakyan, A., Lukianchikov, A., Plath, K., Jacobsen, S. E. and Clark, A. T.** (2016). Naive human pluripotent cells feature a methylation landscape devoid of blastocyst or germline memory. *Cell Stem Cell* **18**, 323-329. doi:10.1016/j.stem.2016.01.019
- Pastor, W. A., Liu, W., Chen, D., Ho, J., Kim, R., Hunt, T. J., Lukianchikov, A., Liu, X., Polo, J. M., Jacobsen, S. E. et al.** (2018). TFAP2C regulates transcription in human naive pluripotency by opening enhancers. *Nat. Cell Biol.* **20**, 553-564. doi:10.1038/s41556-018-0089-0
- Petropoulos, S., Edsgård, D., Reinius, B., Deng, Q., Panula, S. P., Codeluppi, S., Reyes, A. P., Linnarsson, S., Sandberg, R. and Lanner, F.** (2016). Single-cell RNA-seq reveals lineage and X chromosome dynamics in human preimplantation embryos. *Cell* **165**, 1012-1026. doi:10.1016/j.cell.2016.03.023
- Piper, D. E., Batchelor, A. H., Chang, C.-P., Cleary, M. L. and Wolberger, C.** (1999). Structure of a HoxB1-Pbx1 heterodimer bound to DNA. *Cell* **96**, 587-597. doi:10.1016/S0092-8674(00)80662-5
- Pozzi, L., Hodgson, J. A., Burrell, A. S., Sterner, K. N., Raam, R. L. and Disotell, T. R.** (2014). Primate phylogenetic relationships and divergence dates inferred from complete mitochondrial genomes. *Mol. Phylogenet. Evol.* **75**, 165-183. doi:10.1016/j.ympev.2014.02.023
- Riesenberg, S., Chintalapati, M., Macak, D., Kanis, P., Maricic, T. and Pääbo, S.** (2019). Simultaneous precise editing of multiple genes in human cells. *Nucleic Acids Res.* **47**, e116. doi:10.1093/nar/gkz669
- Rostovskaya, M.** (2022). Maintenance of human naïve pluripotent stem cells. *Methods Mol. Biol.* **2416**, 73-90. doi:10.1007/978-1-0716-1908-7_6
- Rostovskaya, M., Stirparo, G. G. and Smith, A.** (2019). Capacitation of human naive pluripotent stem cells for multi-lineage differentiation. *Development* **146**, dev172916. doi:10.1242/dev.172916
- Rostovskaya, M., Andrews, S., Reik, W. and Rugg-Gunn, P. J.** (2022). Amniogenesis occurs in two independent waves in primates. *Cell Stem Cell* **29**, 744-759.e6. doi:10.1016/j.stem.2022.03.014
- Rugg-Gunn, P. J.** (2022). Induction of human naive pluripotency using chemical resetting. *Methods Mol. Biol.* **2416**, 29-37. doi:10.1007/978-1-0716-1908-7_3
- Sayers, E. W., Bolton, E. E., Brister, J. R., Canese, K., Chan, J., Comeau, D. C., Connor, R., Funk, K., Kelly, C., Kim, S. et al.** (2022). Database resources of the national center for biotechnology information. *Nucleic Acids Res.* **50**, D20-D26. doi:10.1093/nar/gkab1112
- Scerbo, P., Markov, G. V., Vivien, C., Kodjabachian, L., Demeneix, B., Coen, L. and Girardot, F.** (2014). On the origin and evolutionary history of NANOG. *PLoS ONE* **9**, e85104. doi:10.1371/journal.pone.0085104
- Shakiba, N., White, C. A., Lipsitz, Y. Y., Yachie-Kinoshita, A., Tonge, P. D., Hussein, S. M. I., Puri, M. C., Elbaz, J., Morrissey-Scout, J., Li, M. et al.** (2015). CD24 tracks divergent pluripotent states in mouse and human cells. *Nat. Commun.* **6**, 7329. doi:10.1038/ncomms8329
- Skarnes, W. C., Pellegrino, E. and McDonough, J. A.** (2019). Improving homology-directed repair efficiency in human stem cells. *Methods* **164-165**, 18-28. doi:10.1016/j.ymeth.2019.06.016
- Sozen, B., Jorgensen, V., Weatherbee, B. A. T., Chen, S., Zhu, M. and Zernicka-Goetz, M.** (2021). Reconstructing aspects of human embryogenesis with pluripotent stem cells. *Nat. Commun.* **12**, 5550. doi:10.1038/s41467-021-25853-4
- Štefková, K., Procházková, J. and Pacherník, J.** (2015). Alkaline phosphatase in stem cells. *Stem Cells Int.* **2015**, 1-11. doi:10.1155/2015/628368
- Stirparo, G. G., Boroviak, T., Guo, G., Nichols, J., Smith, A. and Bertone, P.** (2018). Integrated analysis of single-cell embryo data yields a unified transcriptome signature for the human preimplantation epiblast. *Development* **145**, dev158501. doi:10.1242/dev.158501
- Takashima, Y., Guo, G., Loos, R., Nichols, J., Ficiz, G., Krueger, F., Oxley, D., Santos, F., Clarke, J., Mansfield, W. et al.** (2014). Resetting transcription factor control circuitry toward ground-state pluripotency in human. *Cell* **158**, 1254-1269. doi:10.1016/j.cell.2014.08.029
- Tamura, K., Dudley, J., Nei, M. and Kumar, S.** (2007). MEGA4: Molecular Evolutionary Genetics Analysis (MEGA) software version 4.0. *Mol. Biol. Evol.* **24**, 1596-1599. doi:10.1093/molbev/msm092
- Theunissen, T. W., Costa, Y., Radzishuevskaya, A., van Oosten, A. L., Laval, F., Pain, B., Castro, L. F. C. and Silva, J. C. R.** (2011). Reprogramming capacity of Nanog is functionally conserved in vertebrates and resides in a unique homeodomain. *Development* **138**, 4853-4865. doi:10.1242/dev.068775
- Theunissen, T. W., Powell, B. E., Wang, H., Mitalipova, M., Faddah, D. A., Reddy, J., Fan, Z. P., Maetzel, D., Ganz, K., Shi, L. et al.** (2014). Systematic

- identification of culture conditions for induction and maintenance of naive human pluripotency. *Cell Stem Cell* **15**, 524–526. doi:10.1016/j.stem.2014.09.003
- Theunissen, T. W., Friedli, M., He, Y., Planet, E., O’Neil, R. C., Markoulaki, S., Pontis, J., Wang, H., Iouranova, A., Imbeault, M. et al. (2016). Molecular criteria for defining the naive human pluripotent state. *Cell Stem Cell* **19**, 502–515. doi:10.1016/j.stem.2016.06.011
- Thomson, J. A., Itskovitz-Eldor, J., Shapiro, S. S., Waknitz, M. A., Swiergiel, J. J., Marshall, V. S. and Jones, J. M. (1998). Embryonic stem cell lines derived from human blastocysts. *Science* **282**, 1145–1147. doi:10.1126/science.282.5391.1145
- Untergasser, A., Cutcutache, I., Koressaar, T., Ye, J., Faircloth, B. C., Remm, M. and Rozen, S. G. (2012). Primer3—new capabilities and interfaces. *Nucleic Acids Res.* **40**, e115. doi:10.1093/nar/gks596
- Vallier, L., Mendjan, S., Brown, S., Chng, Z., Teo, A., Smithers, L. E., Trotter, M. W. B., Cho, C. H.-H., Martinez, A., Rugg-Gunn, P. et al. (2009). Activin/Nodal signalling maintains pluripotency by controlling Nanog expression. *Development* **136**, 1339–1349. doi:10.1242/dev.033951
- Villar, D., Berthelot, C., Aldridge, S., Rayner, T. F., Lukk, M., Pignatelli, M., Park, T. J., Deaville, R., Erichsen, J. T., Jasinska, A. J. et al. (2015). Enhancer evolution across 20 mammalian species. *Cell* **160**, 554–566. doi:10.1016/j.cell.2015.01.006
- Wang, W., Lin, C., Lu, D., Ning, Z., Cox, T., Melvin, D., Wang, X., Bradley, A. and Liu, P. (2008). Chromosomal transposition of *PiggyBac* in mouse embryonic stem cells. *Proc. Natl Acad. Sci. USA* **105**, 9290–9295. doi:10.1073/pnas.0801017105
- Weiler, S., Gruschus, J. M., Tsao, D. H. H., Yu, L., Wang, L.-H., Nirenberg, M. and Ferretti, J. A. (1998). Site-directed Mutations in the vnd/NK-2 Homeodomain. *J. Biol. Chem.* **273**, 10994–11000. doi:10.1074/jbc.273.18.10994
- Wojdyła, K., Collier, A. J., Fabian, C., Nisi, P. S., Biggins, L., Oxley, D. and Rugg-Gunn, P. J. (2020). Cell-surface proteomics identifies differences in signaling and adhesion protein expression between naive and primed human pluripotent stem cells. *Stem Cell Rep.* **14**, 972–988. doi:10.1016/j.stemcr.2020.03.017
- Woltjen, K., Michael, I. P., Mohseni, P., Desai, R., Mileikovsky, M., Hämäläinen, R., Cowling, R., Wang, W., Liu, P., Gertsenstein, M. et al. (2009). piggyBac transposition reprograms fibroblasts to induced pluripotent stem cells. *Nature* **458**, 766–770. doi:10.1038/nature07863
- Xiang, L., Yin, Y., Zheng, Y., Ma, Y., Li, Y., Zhao, Z., Guo, J., Ai, Z., Niu, Y., Duan, K. et al. (2020). A developmental landscape of 3D-cultured human pre-gastrulation embryos. *Nature* **577**, 537–542. doi:10.1038/s41586-019-1875-y
- Yanagida, A., Spindlow, D., Nichols, J., Dattani, A., Smith, A. and Guo, G. (2021). Naive stem cell blastocyst model captures human embryo lineage segregation. *Cell Stem Cell* **28**, 1016–1022.e4. doi:10.1016/j.stem.2021.04.031
- Yang, Z. and Nielsen, R. (2000). Estimating synonymous and nonsynonymous substitution rates under realistic evolutionary models. *Mol. Biol. Evol.* **17**, 32–43. doi:10.1093/oxfordjournals.molbev.a026236
- Yu, L., Wei, Y., Duan, J., Schmitz, D. A., Sakurai, M., Wang, L., Wang, K., Zhao, S., Hon, G. C. and Wu, J. (2021). Blastocyst-like structures generated from human pluripotent stem cells. *Nature* **591**, 620–626. doi:10.1038/s41586-021-03356-y
- Zaehres, H., William Lensch, M., Daheron, L., Stewart, S. A., Itskovitz-Eldor, J. and Daley, G. Q. (2005). High-efficiency RNA interference in human embryonic stem cells. *Stem Cells* **23**, 299–305. doi:10.1634/stemcells.2004-0252
- Zetsche, B., Gootenberg, J. S., Abudayyeh, O. O., Slaymaker, I. M., Makarova, K. S., Essletzbichler, P., Volz, S. E., Joung, J., van der Oost, J., Regev, A. et al. (2015). Cpf1 is a single RNA-guided endonuclease of a class 2 CRISPR-Cas system. *Cell* **163**, 759–771. doi:10.1016/j.cell.2015.09.038