

UC Berkeley

UC Berkeley Previously Published Works

Title

Mammalian genome innovation through transposon domestication

Permalink

<https://escholarship.org/uc/item/7564x6kb>

Journal

Nature Cell Biology, 24(9)

ISSN

1465-7392

Authors

Modzelewski, Andrew J
Gan Chong, Johnny
Wang, Ting
[et al.](#)

Publication Date

2022-09-01

DOI

10.1038/s41556-022-00970-4

Peer reviewed



Published in final edited form as:

Nat Cell Biol. 2022 September ; 24(9): 1332–1340. doi:10.1038/s41556-022-00970-4.

Transposon domestication expands the functional reservoir in mammalian genomes

Andrew Modzelewski^{1,2}, Johnny Gan Chong¹, Ting Wang³, Lin He^{1,#}

¹Division of Cellular and Developmental Biology, MCB Department, University of California, Berkeley, Berkeley, CA 94720, USA

²Department of Biomedical Sciences, School of Veterinary Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA

³Department of Genetics, Edison Family Center for Genome Science and System Biology, McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO 63110, USA

Abstract

Since the discovery of transposons, their sheer abundance in host genomes has puzzled many. While historically viewed as largely harmless “parasitic” DNAs during evolution, transposons are not a mere record of ancient genome invasion. Instead, nearly every element of transposon biology has been integrated into host biology. Here we review how sequences introduced by transposon activities provide raw material for genome innovation and document the distinct evolutionary path of each species.

Introduction

Barbara McClintock’s seminal discovery of transposable elements (TEs) was decades ahead of its time¹. She postulated the existence of TEs and speculated their gene regulatory activity long before experimental validation¹. Likewise, Britten and Davidson’s “gene battery” model, a theoretical framework on how repetitive sequences contribute to coordinated gene regulation², was not appreciated until recently. With numerous genomes deciphered^{3–5}, it becomes evident that TE influence is widespread across metazoan genomes.

Approximately 40% of mammalian genomes originate from TEs^{4,6}, including DNA transposons (1–2%) and retrotransposons (~40%), both hijacking cellular machineries to spread in host genomes. DNA transposons employ a “cut and paste” mechanism to integrate into the host genome, while retrotransposons use a “copy and paste” strategy for expansion⁷. In recent evolutionary history, retrotransposon domestication is more frequently observed in mammals than that of DNA transposons. Thus, our review focuses on retrotransposons and

#Correspondence to: lhe@berkeley.edu.

Declaration of interests

The authors declare no competing interests.

their roles in genome architecture and innovation. Readers can refer to several recent reviews for DNA transposons^{8–10}.

Retrotransposons are categorized into two groups: long terminal repeat (LTR) and non-LTR retrotransposons (LINEs and SINEs). LTR retrotransposons contain two identical LTRs, flanking an internal protein-coding region; they frequently undergo homologous recombination to generate solo-LTRs. Among the non-LTR retrotransposons, LINEs encode proteins for retrotransposition, while the non-autonomous, non-coding SINE elements exploit LINE-encoded proteins for retrotransposition¹¹.

Given the potential danger associated with rampant transposition, TE abundance in mammals is counterintuitive¹². TEs and their hosts undergo a constant, on-going arms race. TE's ability to colonize, replicate and spread in host genomes is countered by host's surveillance. Most mammalian TEs have been inactivated via degenerative mutations and/or transcriptional/post-transcriptional silencing. Yet occasionally, TE-host interactions, which initially serve a selfish purpose in TE life cycle, can be repurposed for developmental/physiological host functions (Fig. 1). TE fragments could rewire proximal host gene expression by acting as alternative enhancers, promoters, splicing donors/acceptors and polyadenylation signals (Fig. 1). TE elements that encode proteins and/or non-coding RNAs could be domesticated, contributing neogenes to the host for novel biological functions (Fig.1). An intricate balance is struck between selfish TE properties and domesticated TE functionalities. While host genomes are exposed to risks imposed by TE invasion, they gain opportunities for genome innovation that expands gene regulatory modality, enriches transcript diversity, and diversifies functional reservoirs¹³.

Here, we review the roles of TEs in mammalian development, physiology and evolution, with a focus on *in vivo* functional characterization of specific TE elements, as well as the key challenges and opportunities in the field.

Transposons as a functional reservoir of gene regulatory networks

TE-host interactions that mediate TE transcription, splicing, and translational regulation are preserved during evolution and wired into host gene regulatory networks. When proximal to host genes, specific TEs can serve as cell type-specific gene regulatory sequences^{14–16}(Fig. 2), often conferring species-specific gene regulation, and ultimately, species-specific biological readouts (Fig. 3).

TE-derived sequences are prevalent in/near protein-coding genes. 18.4% mouse and 27.4% human Refseq annotations have at least one isoform harboring a TE-derived sequence in its untranslated regions (UTRs)¹⁷; 37% mouse and 45% human enhancers are predicted to be TE-derived¹⁸. The domesticated TEs as gene regulatory elements confer several distinct mechanisms of gene regulation (Fig. 3). Species-specific TEs yield diversification of gene regulation among species, or through convergent evolution, mediate similar gene regulation in different hosts. Additionally, homologous TE loci provide similar/identical gene regulatory sequences to a cohort of host genes, achieving coordinated gene regulation (Fig. 3). These mechanisms greatly enrich host gene regulatory networks.

Successful TE domestication as gene regulatory sequences depends on its gene regulatory capacity, integration sites, and selective evolutionary advantages. In mammals, functional characterizations of TE-dependent gene regulation were often described in germ cells^{19,20} and preimplantation embryos²¹, which are characterized by potent TE induction due to extensive epigenetic reprogramming. TE-mediated gene regulation is also observed in other developmental systems that lack strong, global TE induction, including neurological, hematopoietic and immune systems^{22–28}.

Transposons as promoters

Transposon promoters have been co-opted to regulate specific host gene isoforms, expand transcript diversity and enrich gene regulatory networks. Using technologies ranging from cDNA library cloning²⁹ to RNA-seq³⁰, hundreds of TE promoters have been identified, generating numerous alternative gene isoforms with distinct expression dynamics and/or altered open reading frames (ORFs). The *in vivo* importance of a TE promoter was first revealed by a mouse-specific MTC promoter (an LTR retrotransposon), which drives an oocyte-specific, N-terminally truncated Dicer isoform, Dicer^O. Dicer^O exhibits a greater enzymatic activity than the canonical Dicer, leading to highly efficient RNA interference (RNAi) during oocyte maturation and thus placing RNAi as the central mechanism for post-transcriptional gene/TE silencing in mouse oocytes³¹. Deletion of this MTC element abolishes Dicer^O expression, causing meiotic spindle defects in oocytes, and ultimately, female infertility^{31,32}. In comparison, other mammals lack Dicer^O, and employ the piRNA pathway instead for posttranscriptional silencing³³. These findings reveal the importance of TEs for evolutionary plasticity of species-specific biological processes non-essential for host viability.

TE promoters can also be repurposed for essential mammalian developmental functions³⁴. In preimplantation embryos, a mouse-specific MT2B2 promoter (an LTR retrotransposon) drives transient, yet potent induction of *Cdk2ap1*^N, an N-terminally truncated *Cdk2ap1* isoform³⁴. Unlike canonical *Cdk2ap1* that suppresses cell proliferation, *Cdk2ap1*^N promotes proliferation. The MT2B2 deletion abolishes *Cdk2ap1*^N in preimplantation embryos, causing reduced cell proliferation, embryonic lethality and impaired implantation³⁴. The essential role of the MT2B2 promoter is surprising, as preimplantation development was presumably normal prior to its integration into the ancestral mouse genome. The persistence of MT2B2 in mouse implicates a selective advantage through increased preimplantation cell proliferation, as induced by *Cdk2ap1*^N. Additional changes likely evolved in the mouse genome to adapt to the MT2B2 integration, ultimately rendering it essential.

Domestication of TE promoters can yield either species-specific or -conserved gene regulation. In the case of *Cdk2ap1*^N³⁴, nearly all mammals harbor a *Cdk2ap1*^N isoform with an evolutionarily conserved ORF³⁴. The mouse-specific MT2B2 promoter drives strong preimplantation induction of *Cdk2ap1*^N³⁴. In pig and cow, a transposon-independent promoter regulates *Cdk2ap1*^N expression, but yields minimal preimplantation expression³⁴. In primates, another transposon-derived promoter, L2a/Charlie4z, generates a modest preimplantation expression of *Cdk2ap1*^N³⁴. The L2a/Charlie4z element is upstream of

Cdk2ap1 in many placental mammal genomes, yet is not employed by mouse, cow or pig to regulate *Cdk2ap1*^N. We speculate that an ancient L2a/Charlie4Z integration yields the *Cdk2ap1*^N isoform in ancestral genomes, and that additional transposon integration and/or L2a/Charlie4z degeneration reprogrammed *Cdk2ap1*^N expression in a species-specific manner.

TE-dependent regulation of *prolactin* expression in endometrium tells a different evolutionary story, in which species-specific transposon promoters mediate a conserved gene expression pattern through convergent evolution^{35,36}. The acquisition of endometrium *prolactin* expression in evolution occurred independently in multiple species by domestication of different transposon promoters, including MER77 in mouse, L1-2a in elephant and MER39 in primate (including human), all of which generate a highly conserved expression pattern^{35,36}.

Transposon-derived promoters provide a powerful mechanism for coordinated gene regulation. As a family of transposons quickly spread through the host genome, transcription factor-binding sites embedded within TE promoters are rapidly propagated. Given their sequence similarities, highly related transposon promoters are often coordinately regulated, achieving co-induction of dozens, if not hundreds, of host gene isoforms^{29,30,34,37,38}(Fig. 3). The capacity of transposon promoters to generate new transcriptional regulation, to create new host gene isoforms, and to rewire gene regulatory network, enables genome innovation, particularly in cell types that are susceptible to transposon induction, such as germ cells, preimplantation embryos and placenta.

Human specific, transposon-dependent gene regulation likely underlies human-specific biology. For example, human and great apes maintain fertility for decades, with male fertility more prolonged than females. Unique to humans and great apes, an ERV9 LTR element was integrated upstream of the *p63* gene ~10–15 million years ago, acting as a testis-specific promoter to drive a p63 isoform with an altered N-terminus³⁹. This ERV9:p63 isoform induces a p53 like pro-apoptotic response to eliminate male germ cells with excessive DNA damage, preserving male fertility in human and great apes³⁹. Germ cell specific ERV9 expression is desirable for its spread as a selfish element in the host, yet unexpectedly, a specific ERV9 was repurposed as a guardian of germ cell genome integrity, providing an evolutionary advantage^{40,41}.

Transposons as enhancers, repressors, insulators, and chromatin boundaries

TEs have contributed extensively to enhancers, repressors, and insulators, as previously reviewed^{38,42,43}. On average, ~20% of cell- or tissue-specific, active chromatin elements in human, mouse, and zebrafish are within TEs^{44–46}. The regulatory potential of TEs as enhancers is largely tied to their sequences for transcription factor binding and host chromatin factor recognition. Many transcription factors, including TP53^{47,48}, Oct4⁴⁹, CTCF⁵⁰, and STAT1³⁸, have a large repertoire of TE-derived binding sites⁵¹. Interestingly, some TEs bear a whole array of transcription factor-binding sites, collectively functioning as regulatory modules⁵². Hence, TEs are a source of genetic material for enhancer evolution.

In addition to generating new gene regulatory events, TEs also provide redundancy and robustness to existing regulatory networks. CTCF is a chromatin factor with many binding sites derived from species-specific TEs⁵⁰. Roughly 20% species-specific chromatin loop anchors and topologically associated domain (TAD) boundaries are CTCF sites encoded by species-specific TEs⁵³ (Fig. 2). Strikingly, ~10% of loop anchors and TAD boundaries functionally conserved between human and mouse are derived from species-specific TEs containing CTCF-binding sites⁵³. This paradox can be explained by TE-mediated CTCF binding site turnover, in which existing CTCF-binding sites can be functionally replaced by a new, redundant CTCF site introduced by a proximal species-specific TE insertion that lacks sequence conservation but instead functionally conserves chromatin organization. Indeed, functional conservation in absence of sequence conservation seems to be the rule rather than the exception in the evolution of gene regulatory networks. It remains to be determined if TE-derived chromatin boundaries primarily contribute to genome innovation or robustness of gene regulation.

Transposon as alternative splicing signals

Transposon-dependent alternative splicing is another widespread phenomenon that contributes to evolutionary innovation on gene structure and function (Fig. 2). In some cases, transposons harbor splicing donors and/or acceptors that mobilize host splicing machineries to generate alternative gene isoforms, enabling the incorporation of transposons as gene exons, contributing to alternative coding sequences and/or UTRs. In other cases, transposon integration into host genes adds unique features of pre-mRNA structure, which alters the canonical splicing pattern to generate new gene isoforms with new biological functions.

Among the best examples is an AluY element that integrated into intron 6 of the *TBXT* gene in the hominoid ancestor genome about 25 million years ago⁵⁴ (Fig. 1). Adjacent to the AluY element is a more ancient AluSx1 element integrated in the reverse orientation, resulting in a hairpin structure within the *TBXT* pre-mRNA that traps exon6 and prevents its incorporation in the mRNA. This generates a hominoid-specific, alternative splicing isoform, *TBXT* *exon6*, whose emergence during primate evolution coincides with tail loss in hominoid. *TBXT* *exon6* expression in mice results in impaired tail development or complete tail loss, supporting that AluY integration is an evolutionary event that caused/contributed to tail loss in hominids⁵⁴. Tail loss likely confers a significant selective advantage, possibly by enhancing locomotion and adopting a non-arboreal lifestyle in primate. Hence, a seemingly random event in the transposon-host interaction may have shaped a major event in hominoid evolution.

Transposon as alternative polyadenylation sites

As most transposons mobilize host Pol II machinery for selfish transcription, transposon-derived polyadenylation signals can generate isoforms with altered 3' UTRs⁵⁵. Since 3'UTRs regulate mRNA stability, translation, localization, trafficking, and protein localization⁵⁵, an altered 3'UTR enables distinct post-transcriptional gene regulation. Alternative polyadenylation can also be coupled with alternative splicing to generate protein isoforms with unique C-termini, ultimately, a different protein function (Fig. 2).

Mouse sex determination is among the best examples illustrating how TE-dependent alternative polyadenylation yields functional diversity⁵⁶. *Sry*, a DNA-binding protein, is an essential factor initiating male sex determination in mammals. *Sry* has been considered a single exon gene for 30 years, until an alternative gene isoform, *Sry-T*, was identified in mice⁵⁶. *Sry-T* is generated by alternative splicing coupled with alternative polyadenylation, in which *Sry* exon1 splices into a transposon-derived, second exon, consisting of an L3 element and 3 tandem LTRs. The *Sry-T* transcription terminates at a polyadenylation signal derived from one of the LTR element. The C-terminal 18aa of the canonical *Sry* isoform encodes a degradation motif, which is replaced by a 15aa, degron-free sequence in the *Sry-T* isoform. This mechanism renders *Sry-T* a more stable isoform, reinforcing male specification in mice⁵⁶. Deletion of this transposon-derived *Sry-T* exon 2 in XY mice abolishes *Sry-T* expression, causing male-to-female sex reversal. Hence, the acquisition of an alternative, TE-derived polyadenylation signal for *Sry* confers a mouse-specific functionality in sex determination.

Transposons as a functional reservoir of non-coding RNAs and proteins

In addition to integrating into the host gene regulatory network, domesticated TEs also generate ncRNAs and/or proteins for the host functional repertoires (Fig.1, 4). As such, TEs are often mutated/truncated, retaining minimal sequences for encoding ncRNAs and/or proteins, and through domestication, aspects of their ancestral functions that support TE-host interactions evolved to regulate unique host cellular processes.

Domestication of transposon-encoded non-coding RNAs

In mammals, many ncRNAs contain TE sequences⁵⁷. Preimplantation-specific LINE1 expression plays an important role in chromatin organization during mouse zygotic genome activation (ZGA). Prolonged transcriptional activation of LINE1 or premature transcriptional silencing of LINE1 in mouse zygotes results in developmental arrest. Surprisingly, this effect is not attributed to LINE1-encoded proteins. Instead, LINE1 ncRNAs regulate the dynamic global chromatin accessibility in early mouse embryos⁵⁸. Hence, TE expression is highly regulated in preimplantation development, rather than a consequence of extensive epigenetic reprogramming. Similarly, the human HERV-H retrotransposons lose protein-coding capacity, but exhibit strong RNA expression from >100 loci in human embryonic stem cells (hESCs), where HERV-H long ncRNAs (lncRNAs) establish or maintain pluripotency⁵⁹. Mechanistically, HERV-H lncRNAs act as a nuclear scaffold for transcription factors, transcriptional machineries, and chromatin modifiers, promoting the expression of proximal host genes to sustain pluripotency⁵⁹. Likewise, LINE1 ncRNAs act as a nuclear scaffold to recruit Nucleolin and Kap1 to silence the Dux/MERVL 2-cell transcriptional program and maintain pluripotency gene network in mouse ESCs⁶⁰.

TE-derived ncRNAs have been associated with cancer resistance by promoting innate immune surveillance (Fig. 1). In blind mole rats, premalignant cells experience a global loss of DNA methylation, which triggers retrotransposon induction, generates cytoplasmic RNA/DNA hybrids, and activates the cGAS-STING pathway to eliminate malignant cells⁶¹. Similarly, treating human cancer cells with DNMT inhibitor, 5-Azacytidine, yields

retrotransposon induction, which generates cytoplasmic double-stranded RNAs (dsRNAs) and triggers the RNA sensing pathway to promote type I Interferon response⁶². In both cases, the pathogenic properties of TE ncRNAs serve as a sensor for disease state, triggering innate immune response to eliminate cells with inappropriate TE induction. It is unclear if such benefit is co-opted by the host, or a side effect of harboring transposons by the host.

Domestication of transposon-encoded proteins

Both DNA transposon- and retrotransposon-encoded proteins are co-opted in mammalian genomes, yet annotated retrotransposon proteins are greater in numbers due to their recent domestication. Ancestral LTR retrotransposons and LINEs express proteins to mediate retrotransposition, most of which undergo deleterious mutations and/or epigenetic silencing. Nevertheless, a subset of LTR retrotransposons, particularly endogenous retroviruses (ERVs), retain protein-coding capacity. Among 19 mammalian species examined, 0.05%–0.15% of ERVs retain protein-coding capacity of retroviral origin⁶³. Since the origin of anciently domesticated transposon-derived proteins may not be easily recognizable, both DNA transposon- and retrotransposon-derived protein-coding genes could be underestimated in numbers.

Some retrotransposons retain the protein-coding capacity of Gag, Pol and Env proteins of the retroviral origin. The domestication of retrotransposon-encoded proteins likely enriches host cellular functions and empowers the host to resist invasion by similar TEs⁶⁴. A recent genome analyses in 700 vertebrate genomes uncovered 177 independent co-option events for retroviral protein-coding genes, with the majority being Gag and Env⁶³. Many of these events are retained for a short evolutionary time frame. Similar functionality of ERV proteins can be repeatedly adopted by different mammals from different ancient retroviruses (Fig. 3)⁶³. Intriguingly, some protein-coding retrotransposons evolve into essential genes, supporting that their invasion provides novel ORFs to fulfill new host functions with a selective advantage (Fig. 4).

1) ERV-encoded Gag proteins—Retrotransposon-encoded Gag was once essential for retroviral packaging and budding. Gag contains three key domains: the N-terminal matrix (MA) domain for plasma membrane binding and virion assembly, the central Capsid (CA) domain for viral capsid core formation, and the Nuclear Capsid (NC) domain for viral RNA packaging. Analysis of all annotated human protein-coding genes reveals dozens of Gag-like genes⁶⁵. In addition to annotated cellular genes with a retrotransposon origin, mammalian genomes also harbor ERV loci with partial or complete protein-coding capacity⁶³. Limited functional studies to date suggest that the molecular functions of domesticated Gag-like proteins all have their roots in those of viral Gag in the retrovirus life cycle⁶⁵.

Arc: Arc (activity-regulated cytoskeleton-associated protein) is a key regulator of synaptic plasticity, long-term learning and memory consolidation. Arc originates from the Ty3/gypsy family Gag gene^{66–68}. Arc self-assembles into a virus-like capsid that encapsulates its own mRNA in extracellular vesicles that are released from active synapses. This mechanism transfers Arc mRNAs into the dendrites of neighboring neurons for localized translation^{67,68} (Fig.4). Despite analogous functions, cellular Arc and *bona fide* retroviral Gag exhibit

mechanistic differences. Arc originates from truncated retrotransposons lacking Env, and therefore relies on a different mechanism of uptake⁶⁹. Unlike retroviral Gag that binds specifically to the viral RNAs, Arc binds its own mRNA, and associates with other cellular mRNAs with a lower affinity⁶⁷. Arc evolves a synaptic function that is atypical of a retroviral Gag function. The key Arc function in regulating AMPA receptor trafficking and membrane density likely stems from unique interactions between its ancestral Ty3/gypsy retrotransposon and host cellular proteins⁷⁰.

During evolution, Gag proteins from two lineages of Ty3/gypsy retrotransposons were independently domesticated, leading to convergent evolution of Arc genes in both Tetrapods and Diptera phyla of the animal kingdom⁶⁶. While vertebrate Arc in mice, human and rat only contains a predicted MA domain and a CA domain to mediate intercellular mRNA transfer between neurons⁶⁸; insect Arc in *Drosophila* contains MA, CA and NC motifs and mediates mRNA transfer among neuromuscular junctions⁶⁷. In both cases, the 3' UTR of Arc mRNA is necessary and sufficient for binding to the Arc protein. The evolutionary origin of Arc provides important insights into mechanisms governing synaptic function and highlights the potential of ancient Gag derived cellular genes for mRNA trafficking.

Peg10: Paternally expressed 10 (Peg10), derived from a Ty3/gypsy LTR retrotransposon, is an evolutionarily conserved, imprinted gene in all eutherian mammals⁷¹. Peg10 is paternally expressed in placenta where deletion in mice caused lethality at embryonic day 9.5 (E9.5), largely due to impaired placental development⁷¹. Interestingly, Peg10 retains retroviral overlapping ORFs, generating two ORFs from the same transcript⁷². Peg10-ORF1 encodes a Gag-like protein containing the CA and NC domains, while PEG10-ORF1/2 encodes a fusion of Gag and Pol generated by a programmed -1 frameshift during the translation of PEG10-ORF1⁷². This mechanism resembles the translation of Gag-Pol in retroviruses, supporting a *bona fide* retrotransposon origin for Peg10.

Similar to Arc, Peg10 encapsulates its own mRNA to form capsid-like particles that are secreted in budding vesicles⁷³. The ability of Peg10 to encapsulate and transport its own mRNA has been exploited to generate a modular platform for mRNA delivery by fusing the *Peg10* 3'UTR motif to cargo mRNA⁷⁴. Pseudotyped Peg10 virus-like particles encapsulate such chimeric RNAs in extracellular vesicles to mediate efficient intercellular transfer, thus employing endogenous proteins to minimize immunogenicity in nucleic acid therapy⁷⁴, and providing an innovative method to complement existing viral delivery systems.

2) ERV-encoded Env proteins—Retroviral Env proteins bind to cell surface receptors to mediate fusion between host and viral membranes, thus determining tissue tropism for infection. Syncytin proteins, derived from *Env* genes of multiple ancestral ERVs, provide a similar function in placenta by promoting cell-cell fusion of mononucleated cytotrophoblasts to establish the multi-nucleated syncytiotrophoblast layer⁷⁵. Syncytiotrophoblast layers are formed during implantation, and maintained throughout gestation to mediate exchange of nutrient, gas and waste between maternal and fetal blood, and shield the fetus from the maternal immune response. Syncytin-mediated fusion of cytotrophoblasts is essential for syncytiotrophoblast maturation and placenta development in mammals (Fig. 4).

During mammalian evolution, Syncytin emerged through at least nine independent domestication events from distinct, species-specific ERVs⁷⁶. While different mammalian *Syncytin* genes are not conserved in protein sequences due to their distinct retroviral origins, they all exhibit placenta-specific expression, retain fusogenic activity and persist in evolution for extended periods (> 10 million years)⁷⁶. Although many placental mammals, including human, have domesticated an Env protein as syncytin to mediate cell-cell fusion⁷⁷, functional studies have only been performed in mice⁷⁸. Most mammals have one *Syncytin* gene and one syncytiotrophoblast layer, yet mice have two *Syncytin* genes (*Syncytin-A* and *Syncytin-B*) and two syncytiotrophoblast layers (ST-I and ST-II), adding to functional redundancy, complexity and robustness (Fig. 4)⁷⁵. *Syncytin-A* and *Syncytin-B* entered the rodent genome approximately 20 million years ago, regulating the formation of ST-I and ST-II, respectively^{78,79}. Deletion of *Syncytin-A* disrupts cell fusion of the ST-I layer in placenta, causing aberrant cell expansion, apoptosis and impaired fetal vascularization, and ultimately, embryonic lethality⁷⁸. In contrast, *Syncytin-B* null placenta displays impaired cell fusion of the ST-II layer, yet the embryos are viable with only limited late-onset growth defects.

Syncytin-B also exhibits immune suppressive activity, an innate property of retroviral Env proteins, likely conferring maternal-fetal tolerance⁸⁰. It is tempting to speculate that the consecutive retroviral gene capture by the rodent genome provides a biological innovation that generates a multilayered placental structure with functional redundancy. The domestication of Syncytin in the ancestral mammals could be a pivotal event for the emergence of placental mammals. The replacement of the ancestral *Syncytin* gene with a new Env gene in each species likely contributes to a species-specific mechanism for placentation.

ERV-encoded Pol and Gag-Pol proteins—Retroviral Pol protein contains several important domains, including the protease (PR) that self-cleaves the polyprotein, the reverse transcriptase (RT) domain that converts the RNA genome into cDNA, and the integrase (IN) domain that integrates the retrotransposon genome into the host genome. Pol domestication occurs at a much lower frequency compared to that of Gag and Env, possibly due to the difficulty in taming RT activity that renders detrimental effects.

Bioinformatic analyses have identified two evolutionarily conserved Pol genes, the *Gypsy integrase-1 (GIN-1)* gene harboring an integrase domain, and the *Cousin of Gypsy integrase-1 (CGIN-1)* containing an RNase H and an integrase domains^{81,82}. *GIN-1* and *CGIN-1* are evolutionarily conserved in mammals, implicating a potential host function. Another example of Pol domestication is *Peg10*, which encodes both Gag and Gag-Pol⁷². While *Peg10* deletion leads to mid-gestation lethality in mice⁷¹, mutation of its Pol protease motif causes perinatal lethality, with fetal and placental growth defects due to impaired fetal vasculature⁸³. *Peg10* is expressed in the three trophoblast layers, but not the surrounding fetal capillary epithelial cells⁸³. Interestingly, *Peg11*, presumably derived from the same retrotransposons as *Peg10*, is specifically expressed in fetal endothelial cells, but not trophoblasts. *Peg11* contains Gag and Pol regions, and its deficiency in mice leads to impaired fetal capillaries in placenta during mid to late gestation, resembling the phenotype caused by protease-deficient *Peg10*^{83,84}. While the exact molecular basis remains elusive,

the protease activity of *Peg10* in trophoblasts and the Pol-like activity of *Peg11* in fetal endothelial cells act at the fetal-maternal interface to safeguard the development of fetal vasculature.

Non-LTR retrotransposon-encoded proteins—Non-LTR retrotransposons, such as LINEs, have protein-coding capacity, yet their ORFs are domesticated less frequently in mammals. L1TD1 is perhaps the best-known example in humans. L1TD1 originates from a co-opted LINE-1 element that was initially integrated into the common ancestor of eutherian mammals, but subsequently lost or pseudogenized multiple times in some species during mammalian evolution⁸⁵. It has been speculated that L1TD1 confers genome defense against LINE-1 and may have later evolved other functions such as pluripotency maintenance⁸⁵.

DNA transposon-encoded proteins—DNA transposons are less abundant and active in modern mammalian genomes compared to their retrotransposon counterparts, yet their domestications have also shaped important developmental/physiological processes in evolution. In jawed vertebrates, RAG1 and RAG2, the key enzymes for V(D)J recombination essential for humoral immunity, are derived from transposase genes of ancient, eukaryotic *Transib* DNA transposons⁸⁶. Thap1, Thap9 and Thap11 represent a family of Zinc finger transcription factors with a DNA-binding domain homologous to *Drosophila* P-element transposase. Mutations in Thap1 causes DYT6 dystonia in mouse and human⁸⁷; Thap11 deletion causes peri-implantation lethality and defects in the inner cell mass in mice⁸⁸; human THAP9 exhibits an active P-element transposase activity, yet its function is unknown⁸⁹.

Intriguingly, several transposon-host fusion genes are evolved due to exon shuffling, which contain a transposase DNA-binding domain and a host-derived KRAB domain⁹⁰. These KRAB-transposase fusions functionally combine DNA-binding specificity with transcriptional repression to repress expression of specific genes⁹⁰. Thus, transposase capture is a recurrent mechanism for gene evolution, providing not only DNA-binding specificity, but also splicing sites for novel fusions.

Transposon-encoded proteins as evolutionary adaption for host defense

A reoccurring theme in TE domestication is their adaptation to provide host defense against similar pathogens. As divergent as prokaryotes and vertebrates, their key enzymes for genome defense could all be traced back to ancient DNA transposons that had once invaded the host genome⁹¹. In addition to RAG1/RAG2 where ancient transposases are repurposed for humoral immunity⁸⁶, multiple CRISPR-Cas components are likely co-opted from DNA transposons⁶⁴. Cas1, a key component of the class I CRISPR-Cas system, is derived from the transposase of a *Casposons* DNA transposon⁶⁴. Cas9, the key component of the class II CRISPR-Cas system is derived from IscB, an RNA-guided DNA nuclease encoded by the IS200/IS605 family of DNA transposons^{92,93}. This family of transposons also encode TnpB, an endonuclease distantly related to IscB, and a possible ancestral protein for Cas12⁹³. Thus, RNA-guided DNA nucleases encoded by transposons are likely ancestors for key enzymatic components of the CRISPR-Cas system.

Retrotransposons have also been co-opted for host defense against pathogens. Env proteins can act as restriction factors against infection from related retroviruses. The Env of a retrotransposon could block the activity of a related Env receptor in infected host cells, a process termed receptor interference⁹⁴. Another interesting example is the HERV-T Env protein, which directly binds the cell surface receptor, monocarboxylate transporter-1 (MCT-1), to block its activity, hence protecting the cells from additional infection by HERV-T. Domesticated HERV-T Env likely contributed to the extinction of HERV-T that circulated in primate genomes for ~25 million years before going extinct ~8 million years ago⁹⁵.

Env-mediated host defense also occurs in human preimplantation embryos. HERVK is transiently induced at ZGA, followed by the translation of its ORFs and assembly of virus-like particles⁹⁶. HERVK encodes multiple ORFs, including Rec, a homolog of HIV Rev. Rec expression leads to induction of the interferon-induced viral restriction factor IFITM1, thereby triggering an innate antiviral response to protect embryos from repeated infection⁹⁶. Similarly, *Supressyn*, an Env gene originating from an HERV-fb insertion, acts as a potential restriction factor against retroviruses in preimplantation embryos of humans and other hominoids⁹⁷.

Gag can also act as a restriction factor. The mouse Fv1 gene likely originates from an ancient MuERV-L Gag gene given their sequence similarity⁹⁸. Fv1 protects the host from a variety of retroviruses, particularly murine leukemia virus (MLV)⁷⁵. The exact antiviral mechanism of Fv1 is unclear, yet Fv1 is shown to target capsid proteins of exogenous MLV, blocking MLV infection after viral entry but before viral integration and provirus formation^{98,99}. It is intriguing that an MLV-unrelated Gag protein restricts MLV infection, implicating an unexpected interaction between these two retrotransposon Gag genes⁹⁹.

Challenges and opportunities for transposon research

Limited read length of genomic sequencing data, underdeveloped computational tools, and suboptimal TE annotations, all contribute to analytical challenges associated with the repetitive nature of TEs. Many adopt a strategy that relies on uniquely mapped TE reads, thus underestimating TE abundance by ignoring numerous multiply mapped reads^{100–102}. TE functional characterization is also complicated by its repetitiveness. While CRISPR-, TALEN- or RNAi-based technologies could target some TE families if the number of loci is optimal, it is difficult to attribute any phenotypes to a specific locus. Conversely, genetic disruption of a single TE locus is technically feasible yet selecting a single TE locus for functional studies is challenging due to ambiguity in TE mapping. Finally, investigating the evolutionary history of a TE family can be hampered by inaccurate TE annotations, particularly in genomes assembled from short sequencing reads. Renewed efforts to sequence complete mammalian genomes and transcriptomes with long reads will undoubtedly advance the field^{103–105}.

The integration, spreading, fixation/elimination of TEs in a host genome document the unique evolutionary history of that species. Once selfish elements, TEs that are domesticated, co-opted and repurposed during evolution have contributed a substantial amount of raw material for host genome innovation. The modern-day koalas present

a unique experimental system to investigate TE endogenization, TE co-option and TE evolution¹⁰⁶, as they are undergoing genomic colonization by an exogenous retrovirus, KoRV¹⁰⁶, which has begun transitioning into an endogenous retrovirus.

Understanding TE-host interactions will yield powerful strategies for gene delivery, gene manipulation and genome engineering. Gene delivery mediated by DNA TEs has long been harnessed for genetic studies^{107,108}. More recently, components of domesticated retrotransposons, such as Peg10, have been engineered as a gene delivery tool for RNA therapy, utilizing their efficient RNA packaging ability and capacity to infect a variety of host cell types without eliciting immune response⁷⁴. The innate host mechanisms that silence TEs, including RNAi¹⁰⁹ and CRISPR^{110–112}, can be reprogrammed to silence or engineer endogenous host genes for therapeutic purposes. These approaches have created numerous possibilities to treat a spectrum of human diseases.

Altogether, TE domestication reveals the evolutionary history of genes, gene regulation and genome organization, and significantly contributes to the molecular basis for species-specific, phenotypic diversity. TE biology enriches our understanding on disease mechanisms and empowers us with new therapeutic strategies. Friends or foes, our intimate relationship with TEs may have shaped who we are as a species and will likely continue to do so as long as we co-evolve with our TEs.

Acknowledgments

We are grateful to Leslie B. King for editing and proofreading this review. A.J.M. is supported by NIH (R00HD096108) and the Siebel Stem Cell Institute. T.W. is supported by NIH (R01HG007175, U24ES026699, U01CA200060, U01HG009391, U41HG010972, and U24HG012070). L.H. is a Thomas and Stacey Siebel Distinguished Chair Professor, supported by an HHMI Faculty Scholar award, a Bakar Fellow award, and NIH grants (1R01GM114414, R01CA139067, 1R21OD027053, GRANT12095758, and R01NS120287).

References

1. McClintock B The Origin and Behavior of Mutable Loci in Maize. *Proc Natl Acad Sci U S A* 36, 344 (1950). [PubMed: 15430309]
2. Britten RJ & Davidson EH Gene Regulation for Higher Cells: A Theory. *Science* (1979) 165, 349–357 (1969).
3. Nurk S et al. The complete sequence of a human genome. *bioRxiv* 2021.05.26.445798 (2021) doi:10.1101/2021.05.26.445798.
4. Lander ES et al. Initial sequencing and analysis of the human genome. *Nature* 2001 409:6822 409, 860–921 (2001).
5. Craig Venter J et al. The sequence of the human genome. *Science* (1979) 291, 1304–1351 (2001).
6. Waterston RH et al. Initial sequencing and comparative analysis of the mouse genome. *Nature* 2003 420:6915 420, 520–562 (2002).
7. Wicker T et al. A unified classification system for eukaryotic transposable elements. *Nature Reviews Genetics* 2007 8:12 8, 973–982 (2007).
8. Jangam D, Feschotte C & Betrán E Transposable Element Domestication As an Adaptation to Evolutionary Conflicts. *Trends in Genetics* 33, 817–831 (2017). [PubMed: 28844698]
9. Feschotte C & Pritham EJ DNA Transposons and the Evolution of Eukaryotic Genomes. *Annu Rev Genet* 41, 331 (2007). [PubMed: 18076328]
10. Feschotte C Transposable elements and the evolution of regulatory networks. *Nat Rev Genet* 9, 397–405 (2008). [PubMed: 18368054]

11. Levin HL & Moran J. v. Dynamic interactions between transposable elements and their hosts. *Nature Reviews Genetics* 2011 12:9 12, 615–627 (2011).
12. Doolittle WF & Sapienza C Selfish genes, the phenotype paradigm and genome evolution. *Nature* 284, 601–603 (1980). [PubMed: 6245369]
13. Gerdes P, Richardson SR, Mager DL & Faulkner GJ Transposable elements in the mammalian embryo: pioneers surviving through stealth and service. *Genome Biology* 2016 17:1 17, 1–17 (2016).
14. Polak P & Domany E Alu elements contain many binding sites for transcription factors and may play a role in regulation of developmental processes. *BMC Genomics* 7, 133 (2006). [PubMed: 16740159]
15. Lowe CB & Haussler D 29 Mammalian Genomes Reveal Novel Exaptations of Mobile Elements for Likely Regulatory Functions in the Human Genome. *PLoS ONE* 7, e43128 (2012). [PubMed: 22952639]
16. Rowe HM et al. TRIM28 repression of retrotransposon-based enhancers is necessary to preserve transcriptional dynamics in embryonic stem cells. *Genome Res* 23, 452–61 (2013). [PubMed: 23233547]
17. van de Lagemaat LN, Landry J-R, Mager DL & Medstrand P Transposable elements in mammals promote regulatory variation and diversification of genes with specialized functions. *Trends in Genetics* 19, 530–536 (2003). [PubMed: 14550626]
18. Simonti CN, Pavli ev M & Capra JA Transposable Element Exaptation into Regulatory Regions Is Rare, Influenced by Evolutionary Age, and Subject to Pleiotropic Constraints. *Molecular Biology and Evolution* 34, 2856–2869 (2017). [PubMed: 28961735]
19. Miyawaki S et al. The mouse Sry locus harbors a cryptic exon that is essential for male sex determination. *Science* (1979) 370, 121–124 (2020).
20. Flemr M et al. A retrotransposon-driven dicer isoform directs endogenous small interfering RNA production in mouse oocytes. *Cell* 155, 807–16 (2013). [PubMed: 24209619]
21. Modzelewski AJ et al. A mouse-specific retrotransposon drives a conserved Cdk2ap1 isoform essential for development. *Cell* 184, 5541–5558.e22 (2021). [PubMed: 34644528]
22. Senft AD & Macfarlan TS Transposable elements shape the evolution of mammalian development. *Nature Reviews Genetics* 2021 22:11 22, 691–711 (2021).
23. Chuong EB, Elde NC & Feschotte C Regulatory activities of transposable elements: from conflicts to benefits. *Nat Rev Genet* 18, 71 (2017). [PubMed: 27867194]
24. Ito J et al. Systematic identification and characterization of regulatory elements derived from human endogenous retroviruses. *PLOS Genetics* 13, e1006883 (2017). [PubMed: 28700586]
25. Notwell JH, Chung T, Heavner W & Bejerano G A family of transposable elements co-opted into developmental enhancers in the mouse neocortex. *Nature Communications* 6, 6644 (2015).
26. Chuong EB, Elde NC & Feschotte C Regulatory evolution of innate immunity through co-option of endogenous retroviruses. *Science* (1979) 351, 1083–1087 (2016).
27. Ye M et al. Specific subfamilies of transposable elements contribute to different domains of T lymphocyte enhancers. *Proc Natl Acad Sci U S A* 117, 7905–7916 (2020). [PubMed: 32193341]
28. Sakashita A et al. Endogenous retroviruses drive species-specific germline transcriptomes in mammals. *Nature Structural & Molecular Biology* 2020 27:10 27, 967–977 (2020).
29. Peaston AE et al. Retrotransposons Regulate Host Genes in Mouse Oocytes and Preimplantation Embryos. *Developmental Cell* 7, 597–606 (2004). [PubMed: 15469847]
30. Macfarlan TS et al. Embryonic stem cell potency fluctuates with endogenous retrovirus activity. *Nature* 487, 57–63 (2012). [PubMed: 22722858]
31. Flemr M et al. A retrotransposon-driven dicer isoform directs endogenous small interfering RNA production in mouse oocytes. *Cell* 155, 807 (2013). [PubMed: 24209619]
32. Franke V et al. Long terminal repeats power evolution of genes and gene expression programs in mammalian oocytes and zygotes. *Genome Research* 27, 1384–1394 (2017). [PubMed: 28522611]
33. Hasuwa H et al. Production of functional oocytes requires maternally expressed PIWI genes and piRNAs in golden hamsters. *Nature Cell Biology* 2021 23:9 23, 1002–1012 (2021).

34. Modzelewski AJ et al. A mouse-specific retrotransposon drives a conserved Cdk2ap1 isoform essential for development. *Cell* 184, 5541–5558.e22 (2021). [PubMed: 34644528]
35. Gerlo S, Davis JRE, Mager DL & Kooijman R Prolactin in man: A tale of two promoters. *BioEssays* 28, 1051–1055 (2006). [PubMed: 16998840]
36. Emera D et al. Convergent Evolution of Endometrial Prolactin Expression in Primates, Mice, and Elephants Through the Independent Recruitment of Transposable Elements. *Molecular Biology and Evolution* 29, 239–247 (2012). [PubMed: 21813467]
37. Davis MP et al. Transposon-driven transcription is a conserved feature of vertebrate spermatogenesis and transcript evolution. *EMBO Rep* 18, 1231–1247 (2017). [PubMed: 28500258]
38. Chuong EB, Elde NC & Feschotte C Regulatory evolution of innate immunity through co-option of endogenous retroviruses. *Science* (1979) 351, 1083–1087 (2016).
39. Beyer U, Moll-Rocek J, Moll UM & Dobbstein M Endogenous retrovirus drives hitherto unknown proapoptotic p63 isoforms in the male germ line of humans and great apes. *Proc Natl Acad Sci U S A* 108, 3624–3629 (2011). [PubMed: 21300884]
40. Pi W et al. The LTR enhancer of ERV-9 human endogenous retrovirus is active in oocytes and progenitor cells in transgenic zebrafish and humans. *Proc Natl Acad Sci U S A* 101, 805–10 (2004). [PubMed: 14718667]
41. Hu T et al. Long non-coding RNAs transcribed by ERV-9 LTR retrotransposon act in cis to modulate long-range LTR enhancer function. *Nucleic Acids Research* 45, gkx055 (2017).
42. Feschotte C Transposable elements and the evolution of regulatory networks. *Nat Rev Genet* 9, 397–405 (2008). [PubMed: 18368054]
43. Sundaram V & Wysocka J Transposable elements as a potent source of diverse cis-regulatory sequences in mammalian genomes. *Philos Trans R Soc Lond B Biol Sci* 375, (2020).
44. Xie M et al. DNA hypomethylation within specific transposable element families associates with tissue-specific enhancer landscape. *Nat Genet* 45, 836–841 (2013). [PubMed: 23708189]
45. Pehrsson EC, Choudhary MNK, Sundaram V & Wang T The epigenomic landscape of transposable elements across normal human development and anatomy. *Nat Commun* 10, (2019).
46. Miao B et al. Tissue-specific usage of transposable element-derived promoters in mouse development. *Genome Biol* 21, (2020).
47. Bourque G et al. Evolution of the mammalian transcription factor binding repertoire via transposable elements. *Genome Res* 18, 1752–62 (2008). [PubMed: 18682548]
48. Wang T et al. Species-specific endogenous retroviruses shape the transcriptional network of the human tumor suppressor protein p53. *Proc Natl Acad Sci U S A* 104, 18613–18618 (2007). [PubMed: 18003932]
49. Kunarso G et al. Transposable elements have rewired the core regulatory network of human embryonic stem cells. *Nat. Genet* 42, 631–634 (2010). [PubMed: 20526341]
50. Schmidt D et al. Waves of retrotransposon expansion remodel genome organization and CTCF binding in multiple mammalian lineages. *Cell* 148, 335–348 (2012). [PubMed: 22244452]
51. Sundaram V et al. Widespread contribution of transposable elements to the innovation of gene regulatory networks. *Genome Res* 24, 1963–1976 (2014). [PubMed: 25319995]
52. Sundaram V et al. Functional cis-regulatory modules encoded by mouse-specific endogenous retrovirus. *Nature Communications* 8, 14550 (2017).
53. Choudhary MNK et al. Co-opted transposons help perpetuate conserved higher-order chromosomal structures. *Genome Biol* 21, (2020).
54. Xia B et al. The genetic basis of tail-loss evolution in humans and apes 2. doi:10.1101/2021.09.14.460388.
55. Mayr C Regulation by 3'-Untranslated Regions. 10.1146/annurev-genet-120116-024704 51, 171–194 (2017).
56. Miyawaki S et al. The mouse Sry locus harbors a cryptic exon that is essential for male sex determination. *Science* (1979) 370, 121–124 (2020).
57. Kelley D & Rinn J Transposable elements reveal a stem cell-specific class of long noncoding RNAs. *Genome Biol* 13, R107 (2012). [PubMed: 23181609]

58. Jachowicz JW et al. LINE-1 activation after fertilization regulates global chromatin accessibility in the early mouse embryo. *Nature Genetics* 2017 49:10 49, 1502–1510 (2017).
59. Lu X et al. The retrovirus HERVH is a long noncoding RNA required for human embryonic stem cell identity. *Nature Structural & Molecular Biology* 2014 21:4 21, 423–425 (2014).
60. Percharde M et al. A LINE1-Nucleolin Partnership Regulates Early Development and ESC Identity. *Cell* 174, 391–405.e19 (2018). [PubMed: 29937225]
61. Zhao Y et al. Transposon-triggered innate immune response confers cancer resistance to the blind mole rat. *Nature Immunology* 2021 22:10 22, 1219–1230 (2021).
62. Chiappinelli KB et al. Inhibiting DNA Methylation Causes an Interferon Response in Cancer via dsRNA Including Endogenous Retroviruses. *Cell* 162, 974–986 (2015). [PubMed: 26317466]
63. Ueda MT et al. Comprehensive genomic analysis reveals dynamic evolution of endogenous retroviruses that code for retroviral-like protein domains. *Mob DNA* 11, (2020).
64. Jangam D, Feschotte C & Betrán E Transposable Element Domestication As an Adaptation to Evolutionary Conflicts. *Trends Genet* 33, 817–831 (2017). [PubMed: 28844698]
65. Campillos M, Doerks T, Shah PK & Bork P Computational characterization of multiple Gag-like human proteins. *Trends Genet* 22, 585–589 (2006). [PubMed: 16979784]
66. Zhang W et al. Structural basis of arc binding to synaptic proteins: implications for cognitive disease. *Neuron* 86, 490–500 (2015). [PubMed: 25864631]
67. Pastuzyn ED et al. The Neuronal Gene Arc Encodes a Repurposed Retrotransposon Gag Protein that Mediates Intercellular RNA Transfer. *Cell* 172, 275–288.e18 (2018). [PubMed: 29328916]
68. Ashley J et al. Retrovirus-like Gag Protein Arc1 Binds RNA and Traffics across Synaptic Boutons. *Cell* 172, 262–274.e11 (2018). [PubMed: 29328915]
69. Kedrov A. v., Durymanov M & Anokhin K v. The Arc gene: Retroviral heritage in cognitive functions. *Neurosci Biobehav Rev* 99, 275–281 (2019). [PubMed: 30772431]
70. Okuno H et al. Inverse synaptic tagging of inactive synapses via dynamic interaction of Arc/Arg3.1 with CaMKII β . *Cell* 149, 886–898 (2012). [PubMed: 22579289]
71. Ono R et al. Deletion of Peg10, an imprinted gene acquired from a retrotransposon, causes early embryonic lethality. *Nature Genetics* 2006 38:1 38, 101–106 (2005).
72. Clark MB et al. Mammalian Gene PEG10 Expresses Two Reading Frames by High Efficiency –1 Frameshifting in Embryonic-associated Tissues. *Journal of Biological Chemistry* 282, 37359–37369 (2007). [PubMed: 17942406]
73. Abed M et al. The Gag protein PEG10 binds to RNA and regulates trophoblast stem cell lineage specification. *PLoS One* 14, (2019).
74. Segel M et al. Mammalian retrovirus-like protein PEG10 packages its own mRNA and can be pseudotyped for mRNA delivery. *Science* (1979) 373, 882–889 (2021).
75. Sha M et al. Syncytin is a captive retroviral envelope protein involved in human placental morphogenesis. *Nature* 2000 403:6771 403, 785–789 (2000).
76. Esnault C, Cornelis G, Heidmann O & Heidmann T Differential Evolutionary Fate of an Ancestral Primate Endogenous Retrovirus Envelope Gene, the EnvV Syncytin, Captured for a Function in Placentation. *PLOS Genetics* 9, e1003400 (2013). [PubMed: 23555306]
77. Sha M et al. Syncytin is a captive retroviral envelope protein involved in human placental morphogenesis. *Nature* 2000 403:6771 403, 785–789 (2000).
78. Dupressoir A et al. Syncytin-A knockout mice demonstrate the critical role in placentation of a fusogenic, endogenous retrovirus-derived, envelope gene. *Proceedings of the National Academy of Sciences* 106, 12127–12132 (2009).
79. Dupressoir A et al. A pair of co-opted retroviral envelope syncytin genes is required for formation of the two-layered murine placental syncytiotrophoblast. *Proc Natl Acad Sci U S A* 108, (2011).
80. Mangeney M et al. Placental syncytins: Genetic disjunction between the fusogenic and immunosuppressive activity of retroviral envelope proteins. *Proceedings of the National Academy of Sciences* 104, 20534–20539 (2007).
81. Marco A & Marín I CGIN1: A Retroviral Contribution to Mammalian Genomes. *Molecular Biology and Evolution* 26, 2167–2170 (2009). [PubMed: 19561090]

82. Lloréns C & Marín I A Mammalian Gene Evolved from the Integrase Domain of an LTR Retrotransposon. *Molecular Biology and Evolution* 18, 1597–1600 (2001). [PubMed: 11470852]
83. Shiura H et al. PEG10 viral aspartic protease domain is essential for the maintenance of fetal capillary structure in the mouse placenta. *Development* 148, (2021).
84. Kitazawa M, Tamura M, Kaneko-Ishino T & Ishino F Severe damage to the placental fetal capillary network causes mid- to late fetal lethality and reduction in placental size in Peg11/Rtl1 KO mice. *Genes to Cells* 22, 174–188 (2017). [PubMed: 28111885]
85. McLaughlin RN et al. Positive selection and multiple losses of the LINE-1-derived L1TD1 gene in mammals suggest a dual role in genome defense and pluripotency. *PLoS Genet* 10, (2014).
86. Huang S et al. Discovery of an Active RAG Transposon Illuminates the Origins of V(D)J Recombination. *Cell* 166, 102 (2016). [PubMed: 27293192]
87. Ruiz M et al. Abnormalities of motor function, transcription and cerebellar structure in mouse models of THAP1 dystonia. *Hum Mol Genet* 24, 7159–7170 (2015). [PubMed: 26376866]
88. Dejosez M et al. Ronin Is Essential for Embryogenesis and the Pluripotency of Mouse Embryonic Stem Cells. *Cell* 133, 1162–1174 (2008). [PubMed: 18585351]
89. Majumdar S, Singh A & Rio DC The human THAP9 gene encodes an active P-element DNA transposase. *Science* 339, 446–448 (2013). [PubMed: 23349291]
90. Cosby RL et al. Recurrent evolution of vertebrate transcription factors by transposase capture. *Science* (1979) 371, (2021).
91. Koonin E. v. & Krupovic M Evolution of adaptive immunity from transposable elements combined with innate immune systems. *Nature Reviews Genetics* 2014 16:3 16, 184–192 (2014).
92. Kapitonov V. v., Makarova KS & Koonin E v. ISC, a Novel Group of Bacterial and Archaeal DNA Transposons That Encode Cas9 Homologs. *J Bacteriol* 198, 797–807 (2015). [PubMed: 26712934]
93. Altae-Tran H et al. The widespread IS200/IS605 transposon family encodes diverse programmable RNA-guided endonucleases. *Science* (1979) 374, 57–65 (2021).
94. Malfavon-Borja R & Feschotte C Fighting Fire with Fire: Endogenous Retrovirus Envelopes as Restriction Factors. *Journal of Virology* 89, 4047 (2015). [PubMed: 25653437]
95. Blanco-Melo D, Gifford RJ & Bieniasz PD Co-option of an endogenous retrovirus envelope for host defense in hominid ancestors. *Elife* 6, (2017).
96. Grow EJ et al. Intrinsic retroviral reactivation in human preimplantation embryos and pluripotent cells. *Nature* 522, 221 (2015). [PubMed: 25896322]
97. Frank JA et al. Antiviral activity of a human placental protein of retroviral origin. *bioRxiv* 2020.08.23.263665 (2020) doi:10.1101/2020.08.23.263665.
98. Yap MW, Colbeck E, Ellis SA & Stoye JP Evolution of the Retroviral Restriction Gene Fv1: Inhibition of Non-MLV Retroviruses. *PLoS Pathogens* 10, (2014).
99. Horikoshi M et al. Positional cloning of the mouse retrovirus restriction gene Fv1. *Nature* 1996 382:6594 382, 826–829 (1996).
100. Jin Y, Tam OH, Paniagua E & Hammell M TEtranscripts: a package for including transposable elements in differential expression analysis of RNA-seq datasets. *Bioinformatics* 31, 3593–3599 (2015). [PubMed: 26206304]
101. Yang WR, Ardeljan D, Pacyna CN, Payer LM & Burns KH SQUIRE reveals locus-specific regulation of interspersed repeat expression. *Nucleic Acids Research* 47, e27–e27 (2019). [PubMed: 30624635]
102. Bendall ML et al. Telescope: Characterization of the retrotranscriptome by accurate estimation of transposable element expression. *PLOS Computational Biology* 15, e1006453 (2019). [PubMed: 31568525]
103. Miga KH et al. Telomere-to-telomere assembly of a complete human X chromosome. *Nature* 585, 79–84 (2020). [PubMed: 32663838]
104. Miga KH & Wang T The Need for a Human Pangenome Reference Sequence. *Annu Rev Genomics Hum Genet* 22, 81–102 (2021). [PubMed: 33929893]
105. Rhie A et al. Towards complete and error-free genome assemblies of all vertebrate species. *Nature* 592, 737–746 (2021). [PubMed: 33911273]

106. Stoye JP Koala retrovirus: a genome invasion in real time. *Genome Biology* 7, 241 (2006). [PubMed: 17118218]
107. Dupuy AJ, Akagi K, Largaespada DA, Copeland NG & Jenkins NA Mammalian mutagenesis using a highly mobile somatic Sleeping Beauty transposon system. *Nature* 436, 221–226 (2005). [PubMed: 16015321]
108. Ding S et al. Efficient transposition of the piggyBac (PB) transposon in mammalian cells and mice. *Cell* 122, 473–483 (2005). [PubMed: 16096065]
109. Zamore PD, Tuschl T, Sharp PA & Bartel DP RNAi: double-stranded RNA directs the ATP-dependent cleavage of mRNA at 21 to 23 nucleotide intervals. *Cell* 101, 25–33 (2000). [PubMed: 10778853]
110. Jinek M et al. RNA-programmed genome editing in human cells. *Elife* 2013, (2013).
111. Mali P et al. RNA-guided human genome engineering via Cas9. *Science* 339, 823–826 (2013). [PubMed: 23287722]
112. Cong L et al. Multiplex genome engineering using CRISPR/Cas systems. *Science* 339, 819–823 (2013). [PubMed: 23287718]
113. Casacuberta E *Drosophila*: Retrotransposons Making up Telomeres. *Viruses* 9, (2017).
114. Abad JP et al. TAHRE, a Novel Telomeric Retrotransposon from *Drosophila melanogaster*, Reveals the Origin of *Drosophila* Telomeres. *Molecular Biology and Evolution* 21, 1620–1624 (2004). [PubMed: 15175413]
115. Levis RW, Ganesan R, Houtchens K, Tolar LA & Sheen F. miin. Transposons in place of telomeric repeats at a *Drosophila* telomere. *Cell* 75, 1083–1093 (1993). [PubMed: 8261510]
116. Xia B et al. The genetic basis of tail-loss evolution in humans and apes. *bioRxiv* 2021.09.14.460388 (2021) doi:10.1101/2021.09.14.460388.

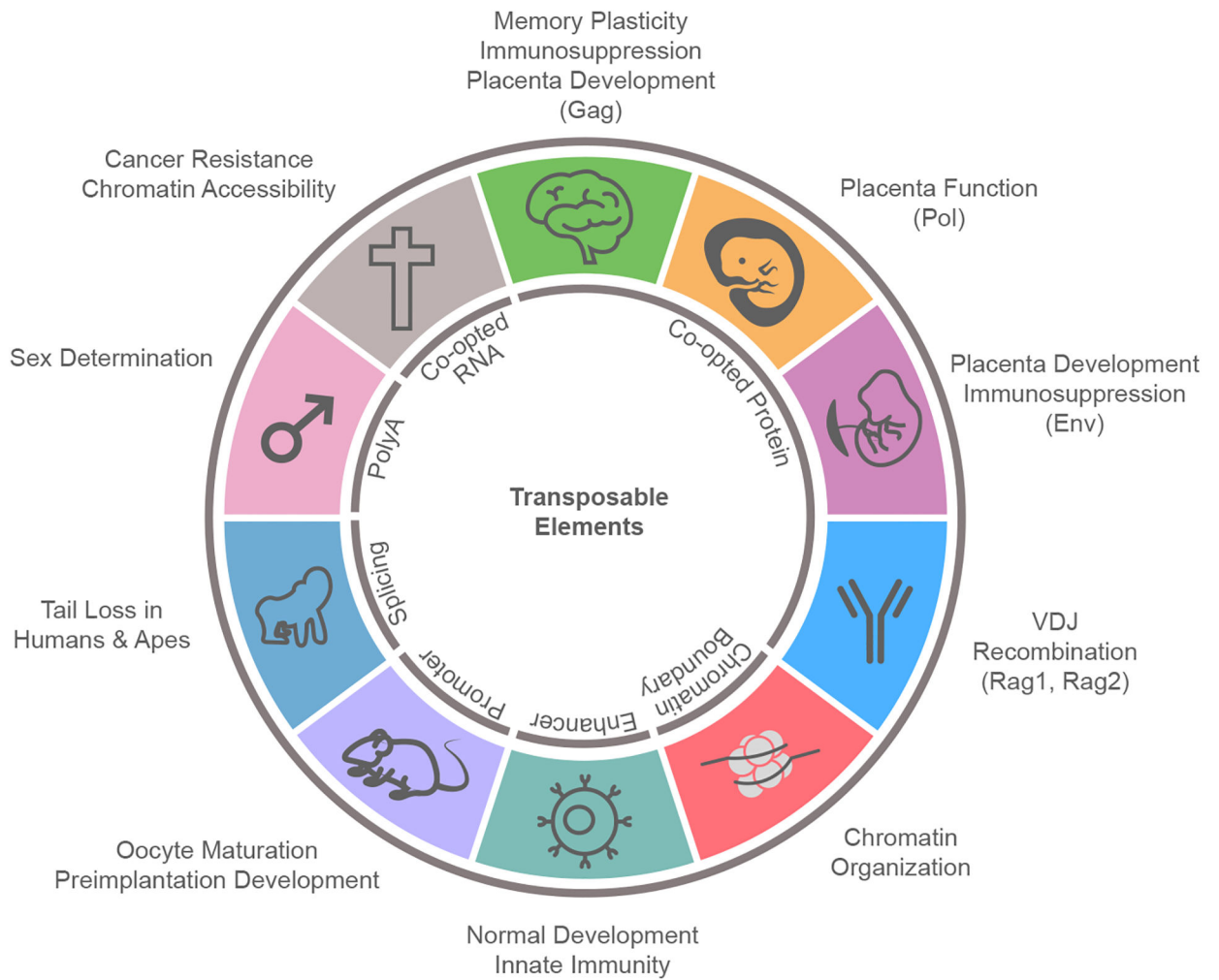


Figure 1. Transposon domestication contributes to host biology.

Transposon domestication provides new mechanisms for host genome innovation in diverse developmental and physiological processes, generating numerous gene regulatory elements, functional ncRNAs and protein-coding genes^{31,32,34,38,42,44,50,52,53,56,58,61,66–68,71,72,75,80,83,84,86,91,96,113–116}. The diagram, while likely representing the tip of an iceberg, summarizes key studies that characterize the *in vivo* validated transposon functions in the host genomes.

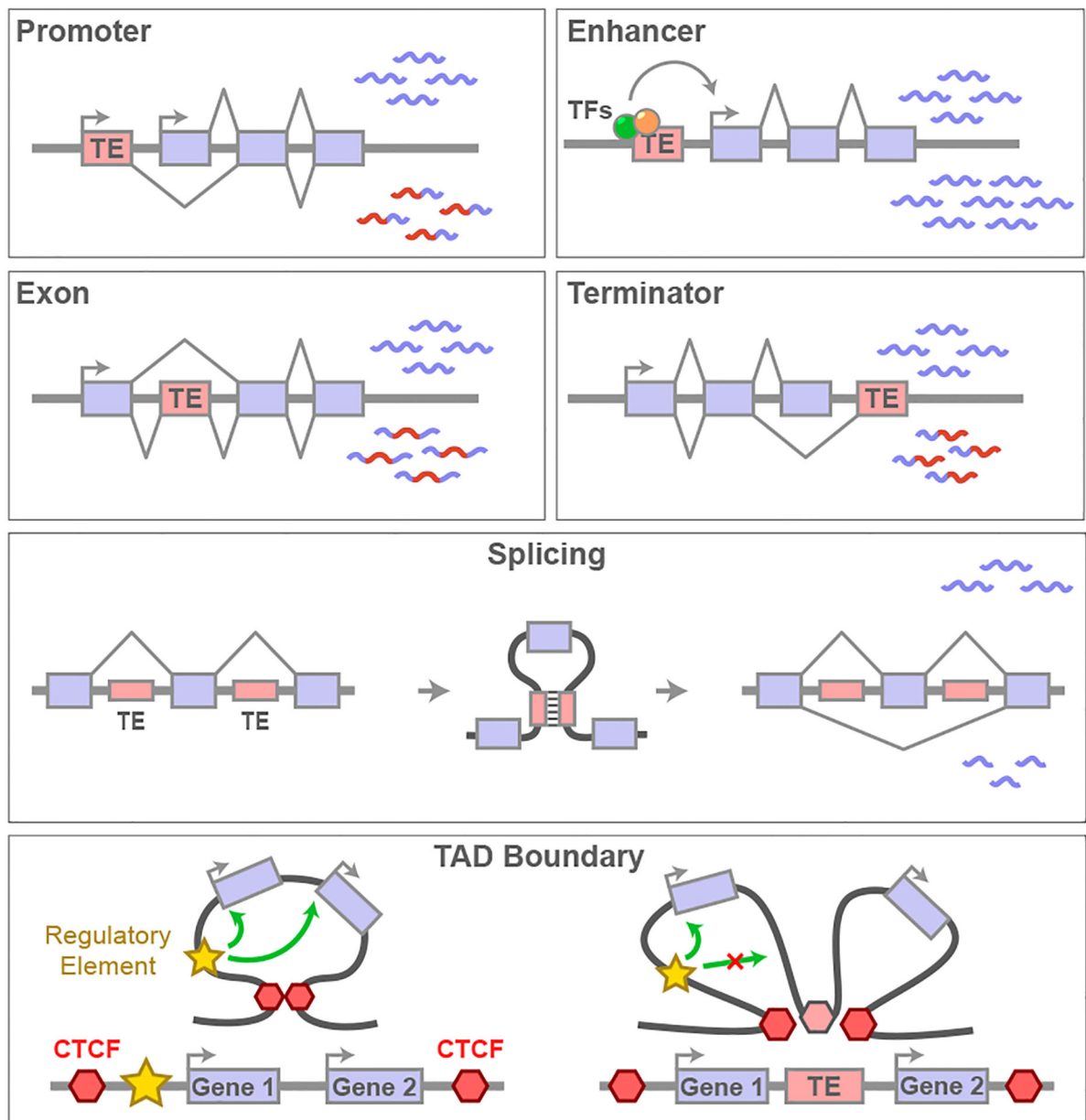


Figure 2. Transposon-derived gene regulatory elements diversify host gene isoforms and enrich expression regulation modality.

Transposon-derived sequences contribute to gene enhancers, promoters, exons, terminators, splicing donors/acceptors, and chromatin boundaries, regulating the structure and expression of proximal host gene isoforms. TE domestication expands gene regulatory modality, enriches transcript diversity, and diversifies functional reservoirs in host genomes. Pink squares, TE elements; blue squares, protein coding exons or protein coding genes; red hexagons, CTCF; yellow star, a gene regulatory element.

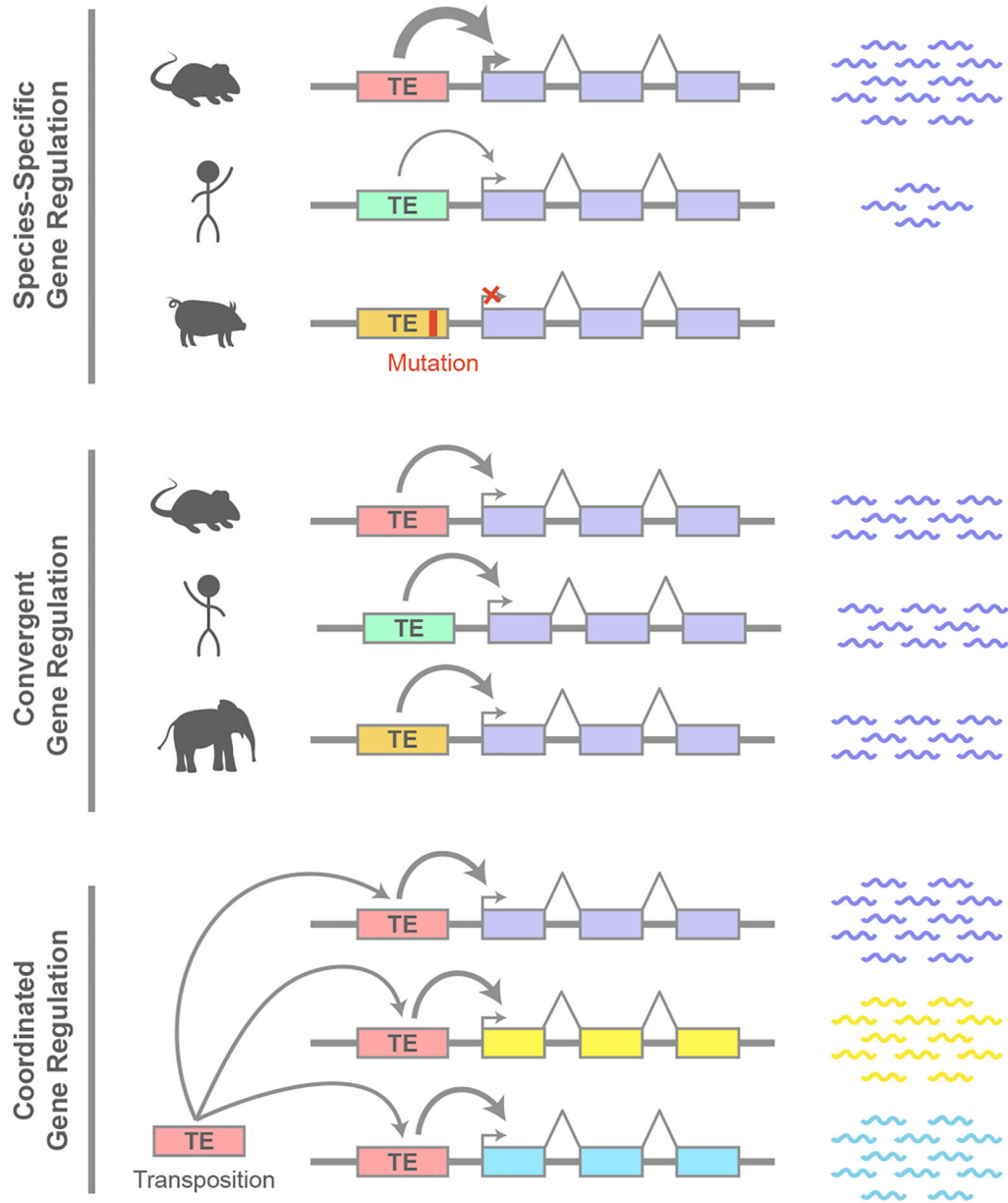


Figure 3. Transposons confer unique modes of cis-gene regulation in host genomes.

Compared to other gene regulatory sequences, transposons have distinct properties, as they are repetitive in nature and frequently species-specific. Species-specific gene regulation occurs when distinct TEs integrate proximal to homologous genes across species, generating a unique expression pattern in each species (top). Convergent gene regulation occurs when distinct TE insertions across species converge on the same regulatory principal to yield nearly identical expression patterns (middle). Coordinated gene regulation occurs when related transposon elements from the same TE family spread in a given host genome and land proximal to a cohort of host genes to coordinate their expression (bottom).

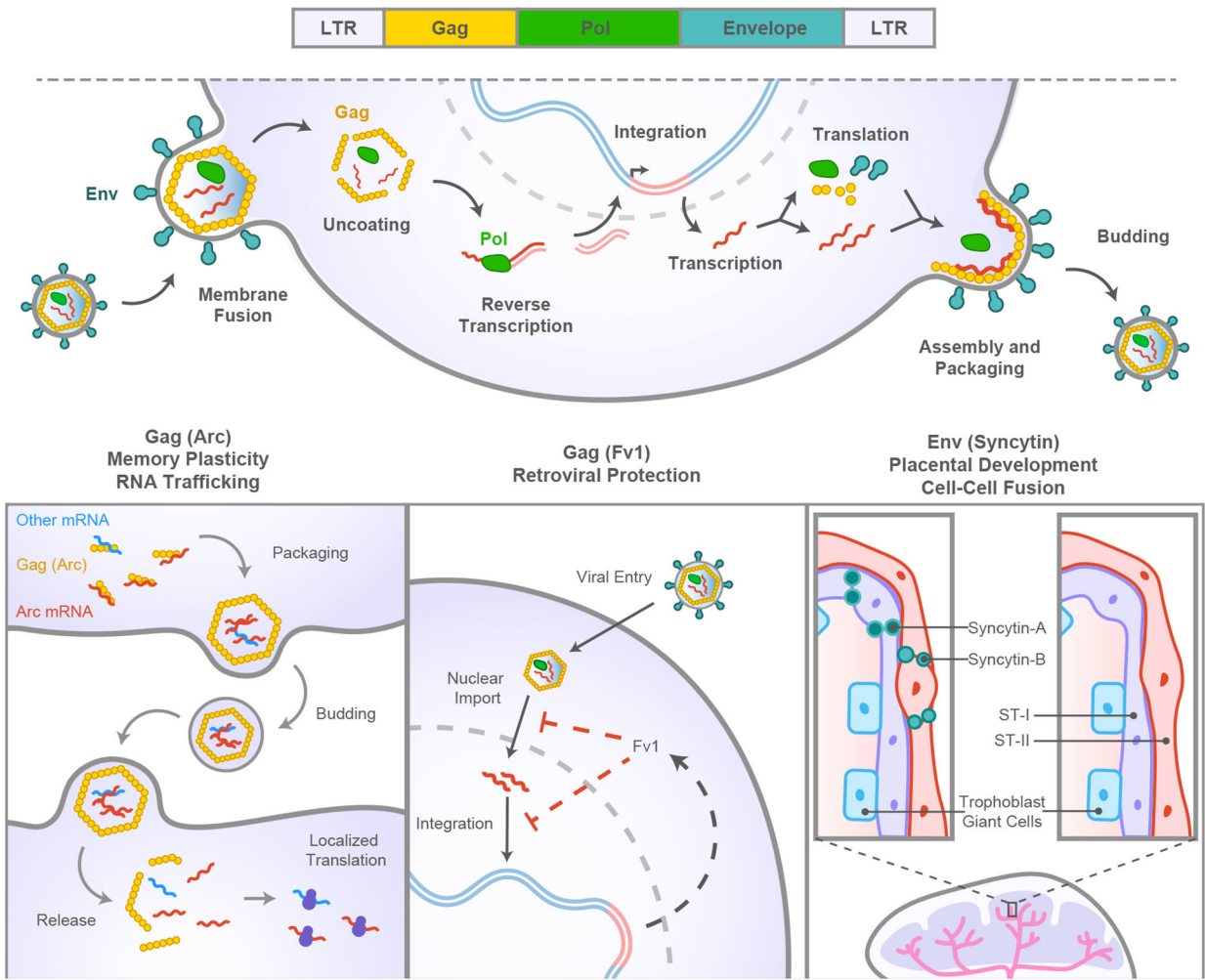


Figure 4: Co-option of transposon-encoded proteins contributes to new host biology.

A diagram illustrates the functional parallel between the retroviral Gag, Pol and Env proteins and their domesticated counterparts encoded by retrotransposons. Retroviral life cycle (top) begins when retroviruses infect the host cells and integrate into the host genome. Subsequently, the host machineries drives the expression of viral Gag, Pol and Env, allowing the retrovirus to mature before released from the host cells. Here, we show examples of domesticated Gag and Env genes (bottom), which are repurposed for neuronal functions, host defense, and placenta development. The remarkable modern innovations conferred by retrotransposon encoded proteins can be traced back to their proviral functions.